

learn network inspire

Game Developers'
Conference

08



February 18-22, 2008
San Francisco

www.gdconf.com



Harnessing the Power of Multiple GPUs

Games Developer Conference 2008
Holger Gruen, holger.gruen@amd.com

Material:

Holger Gruen,

Jon Story, jon.story@amd.com

Ignacio Llamas, illamas@nvidia.com



Agenda

- ④ Why MGPU?
- ④ Driver/Programming Considerations
- ④ Common Pitfalls & Solutions
- ④ Conclusions
- ④ Q&A



Why MGPU?



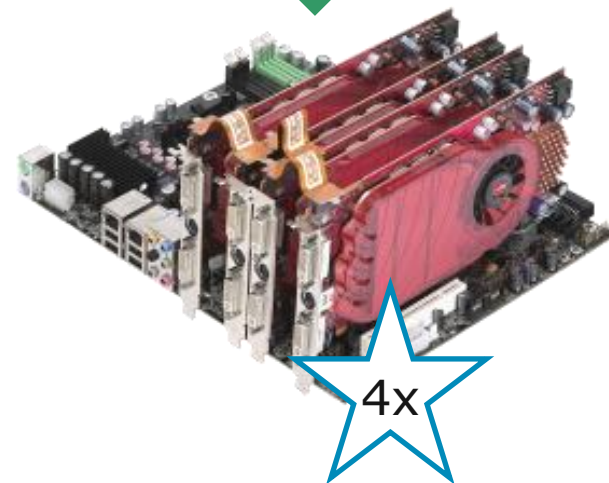
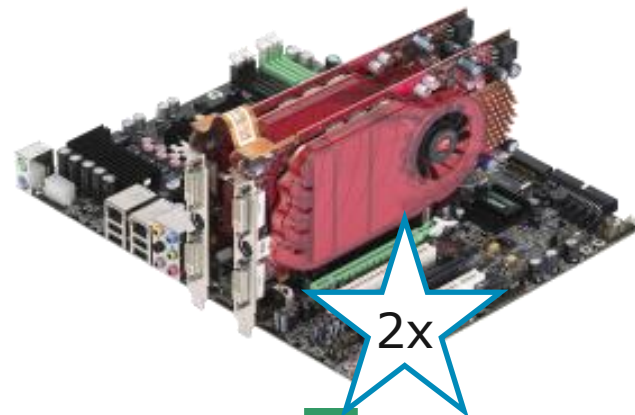
Why MGPU?

- ③ MGPUs can be used to dramatically increase performance and visual quality
 - ③ At higher screen resolutions
 - ③ Especially with increased use of MSAA
- ③ Many applications become GPU limited at higher screen resolutions/ quality settings
 - ③ High resolution monitors => mainstream affordability
- ③ Achieve next generation performance on today's HW
 - ③ Prototype your next engine
 - ③ Speed up your development time



MGPU configurations

- ⊗ Multiple Boards
- ⊗ Multiple GPUs per Board
- ⊗ Hybrid MGPU

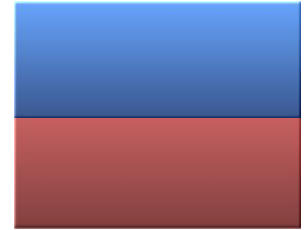




MGPU Rendering Modes

⌚ Split Frame Rendering / Scissor

- ⌚ Screen is divided between GPUs
- ⌚ Dynamic load balancing



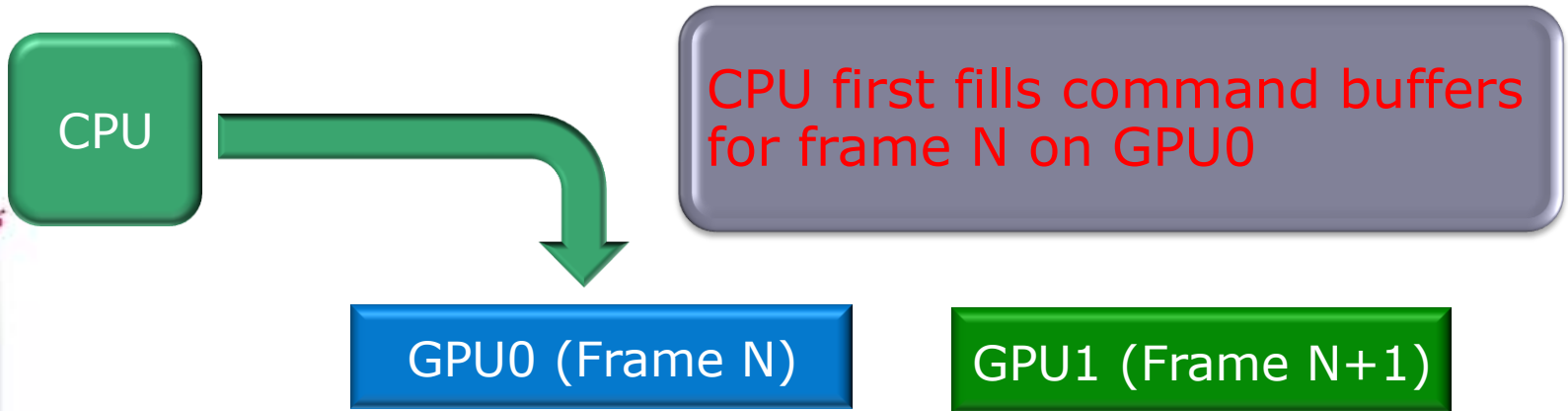
⌚ Alternate Frame Rendering

- ⌚ GPUs take alternate frames
- ⌚ Highest performing mode



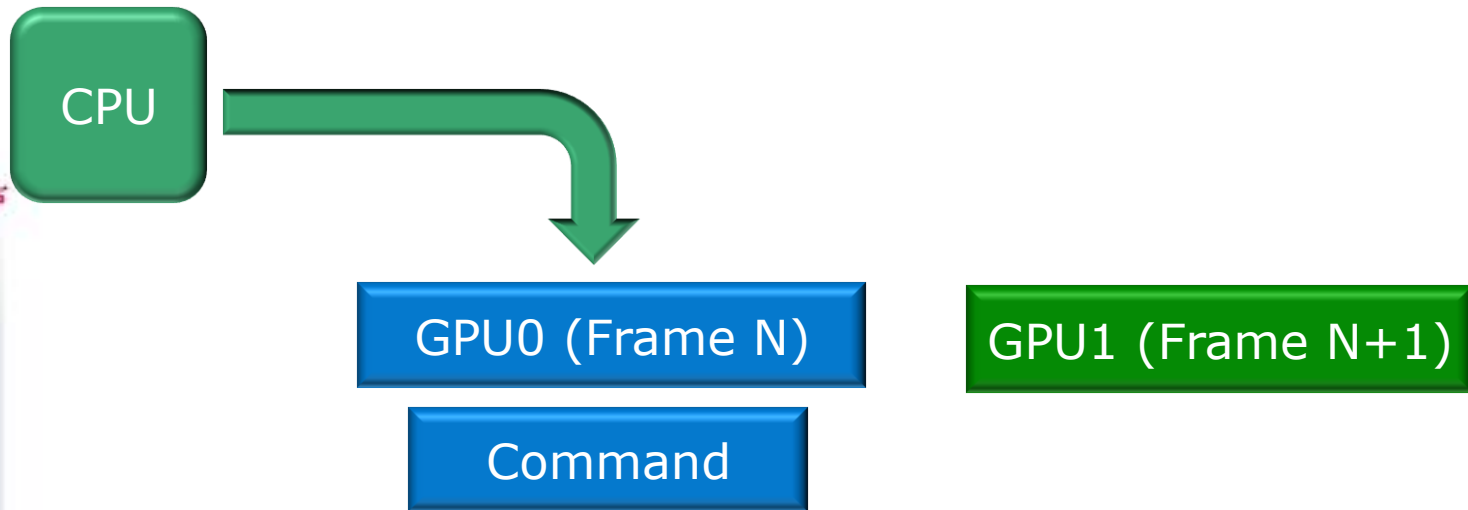


How does AFR Work?



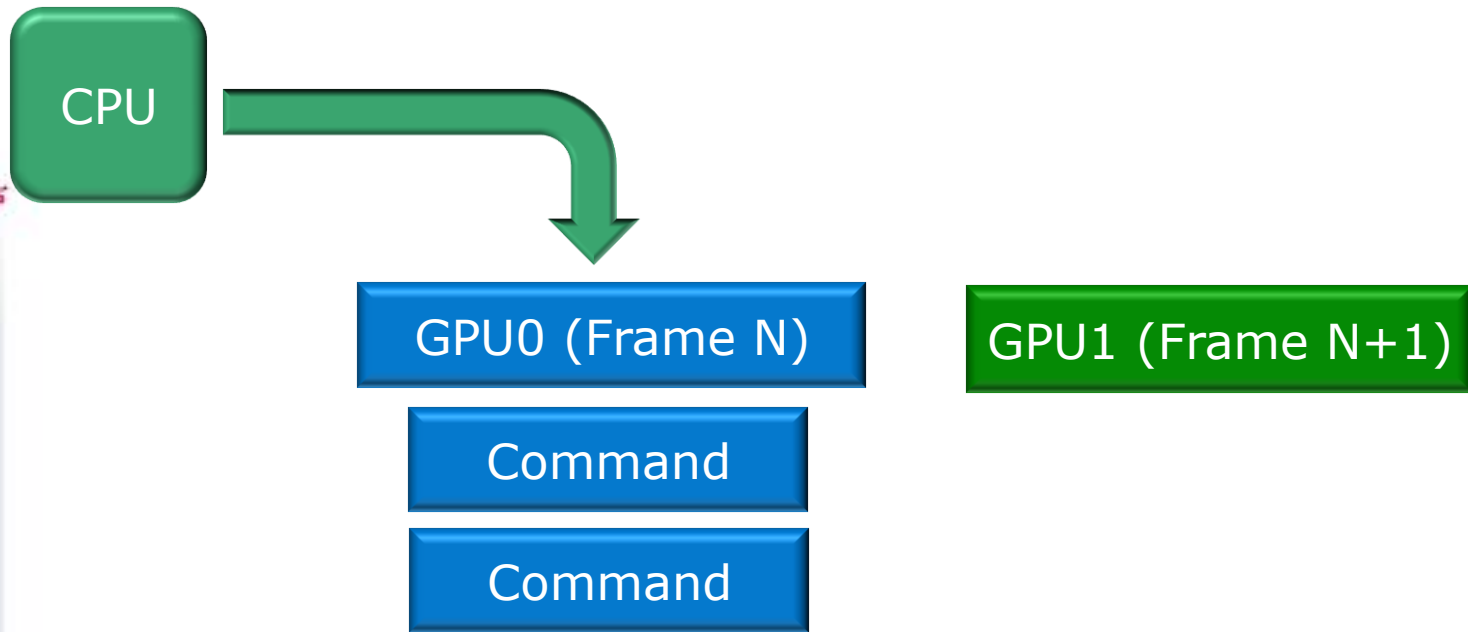


How does AFR Work?



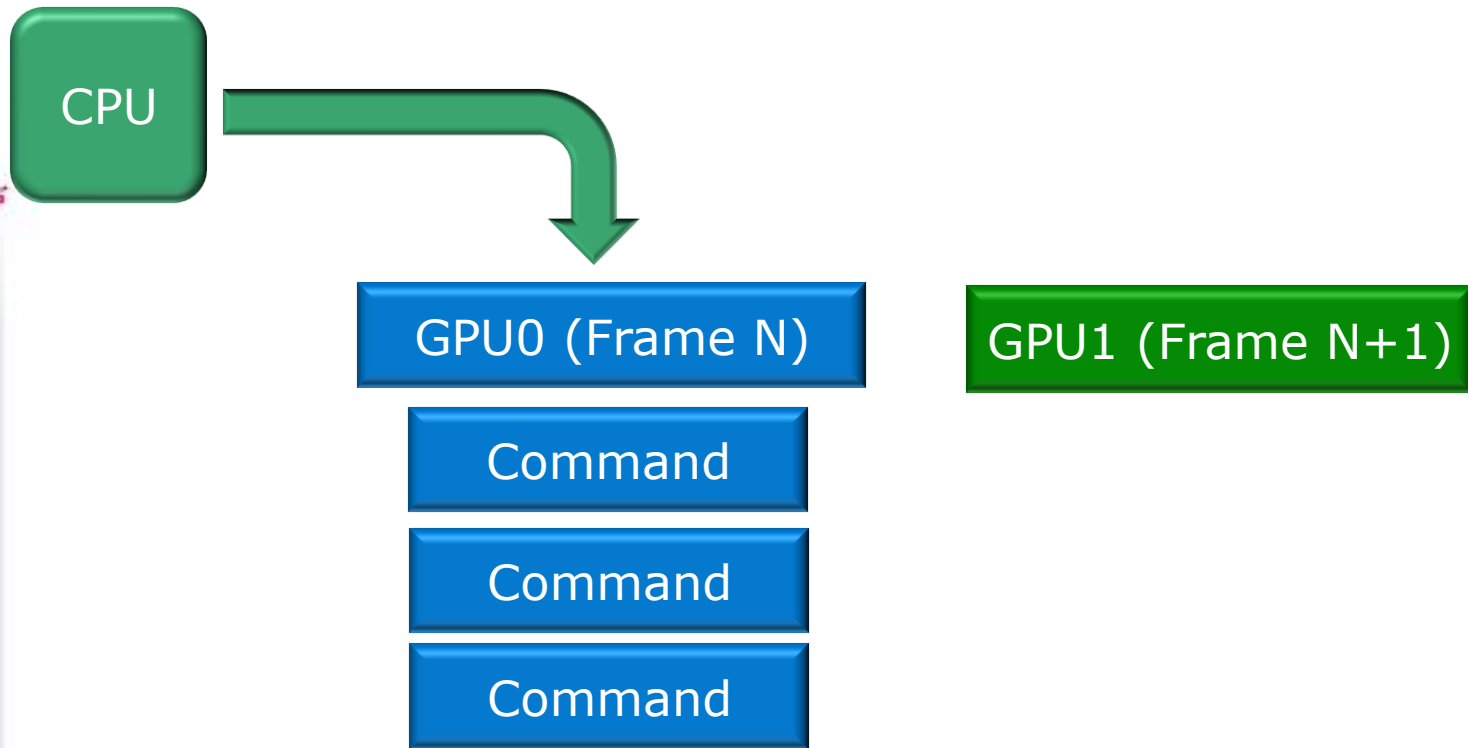


How does AFR Work?



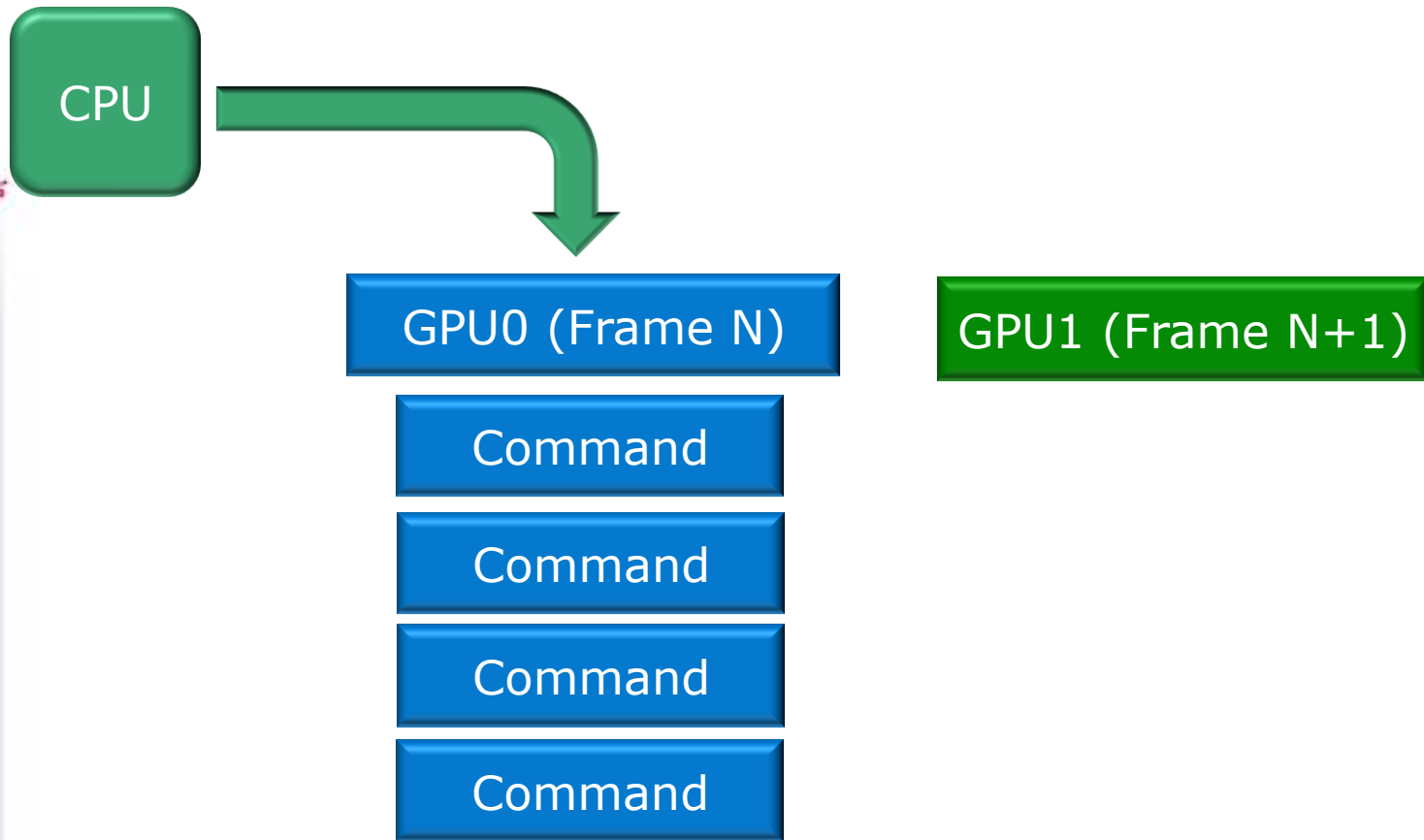


How does AFR Work?



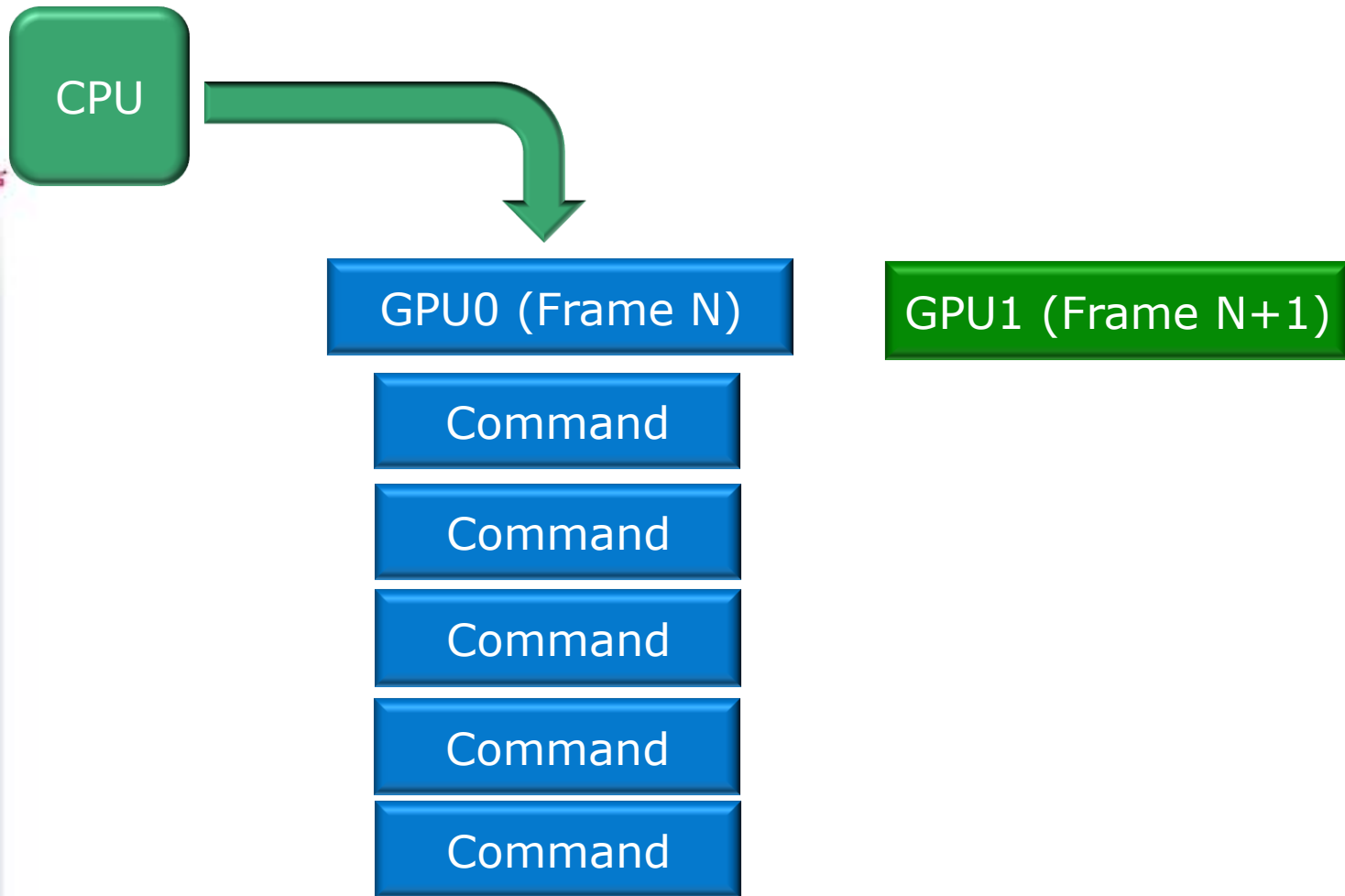


How does AFR Work?



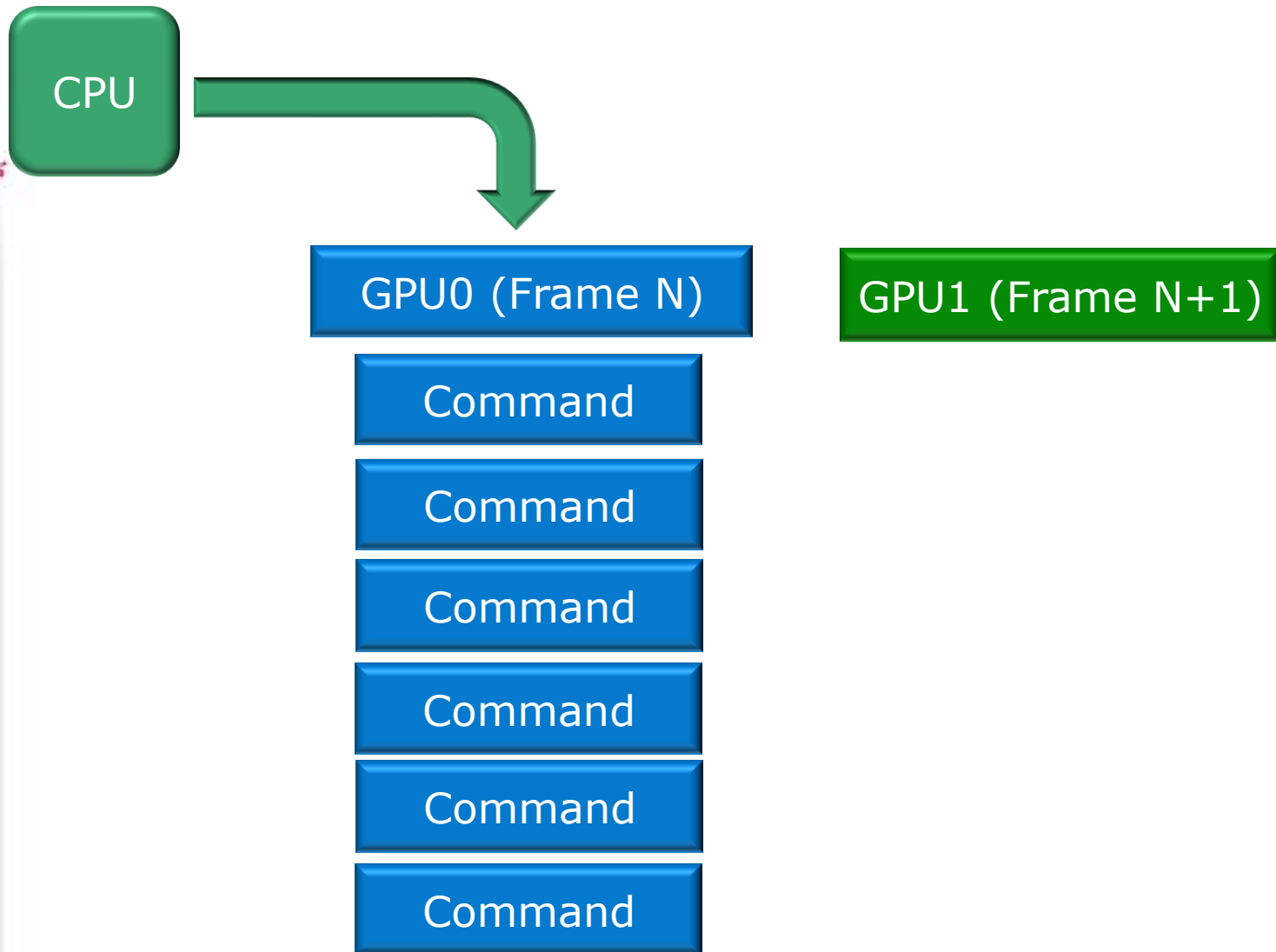


How does AFR Work?



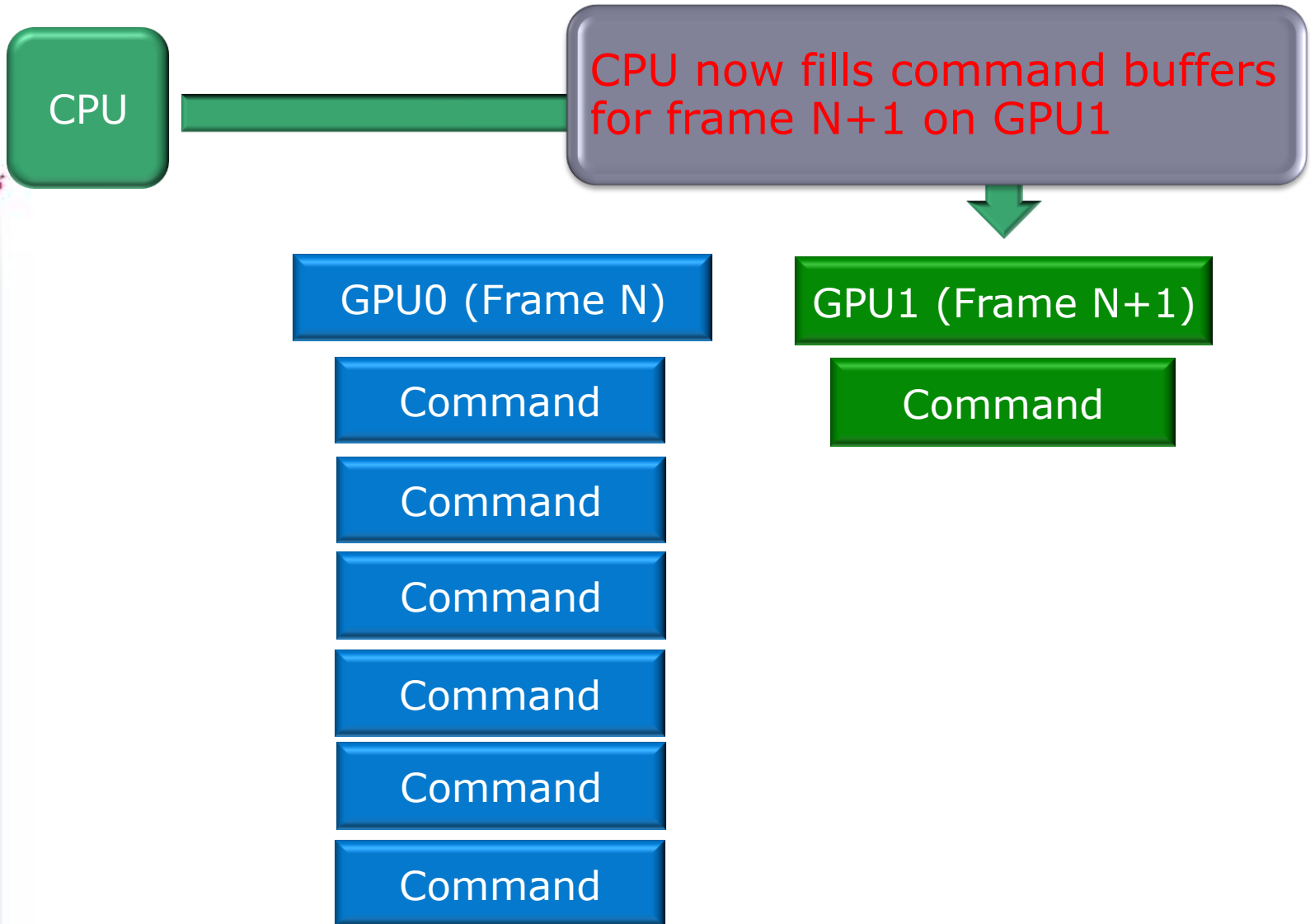


How does AFR Work?



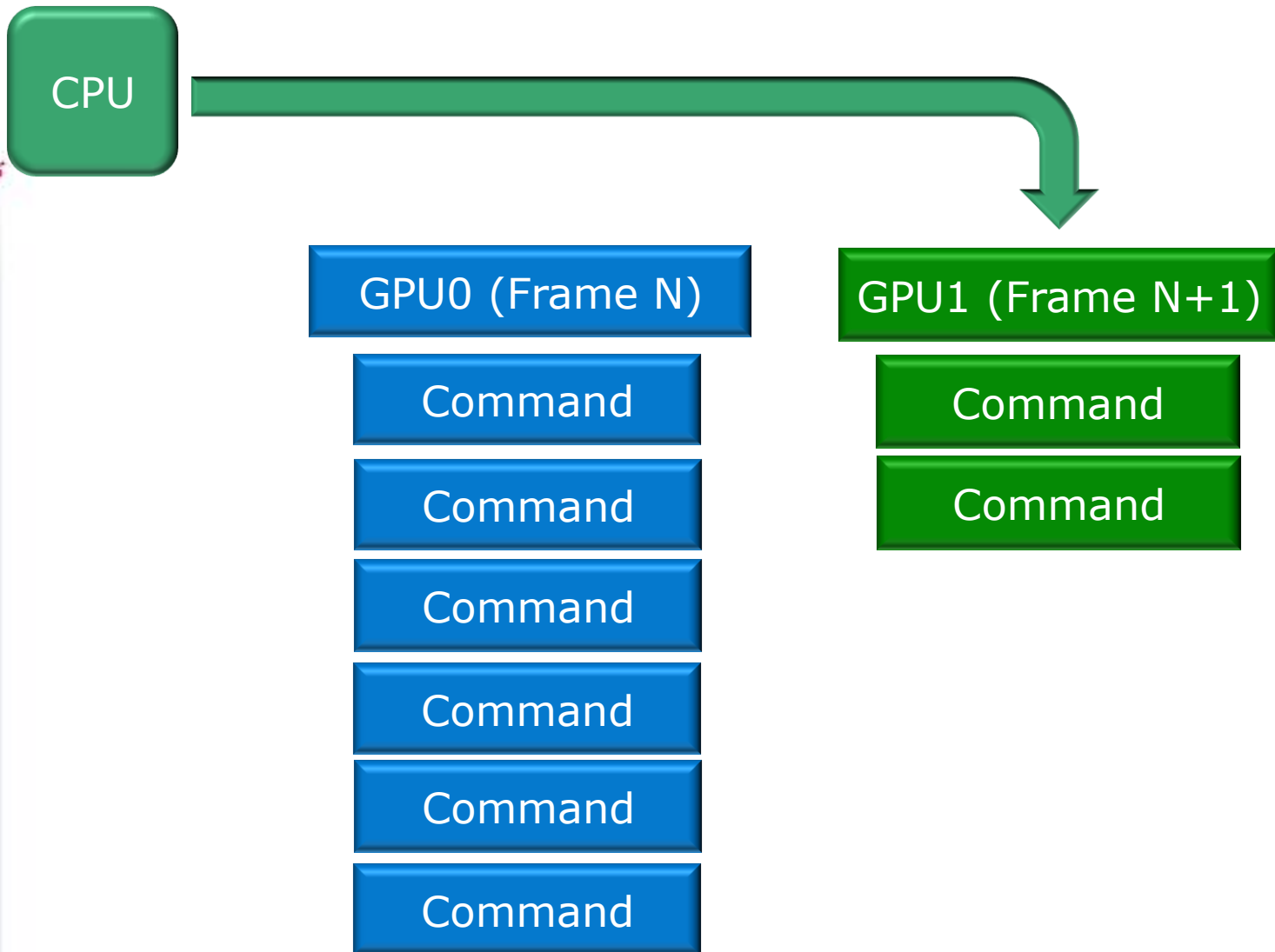


How does AFR Work?



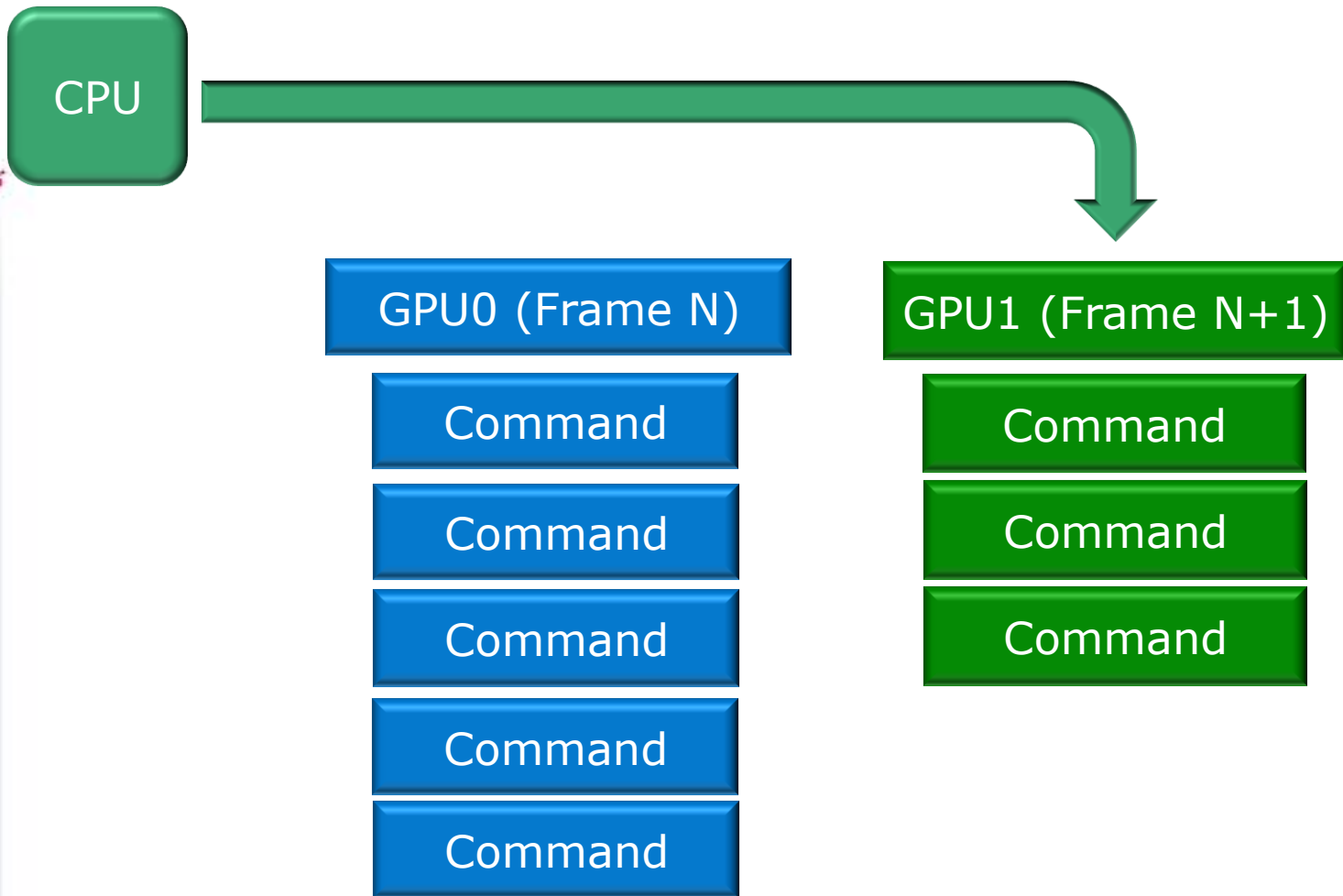


How does AFR Work?



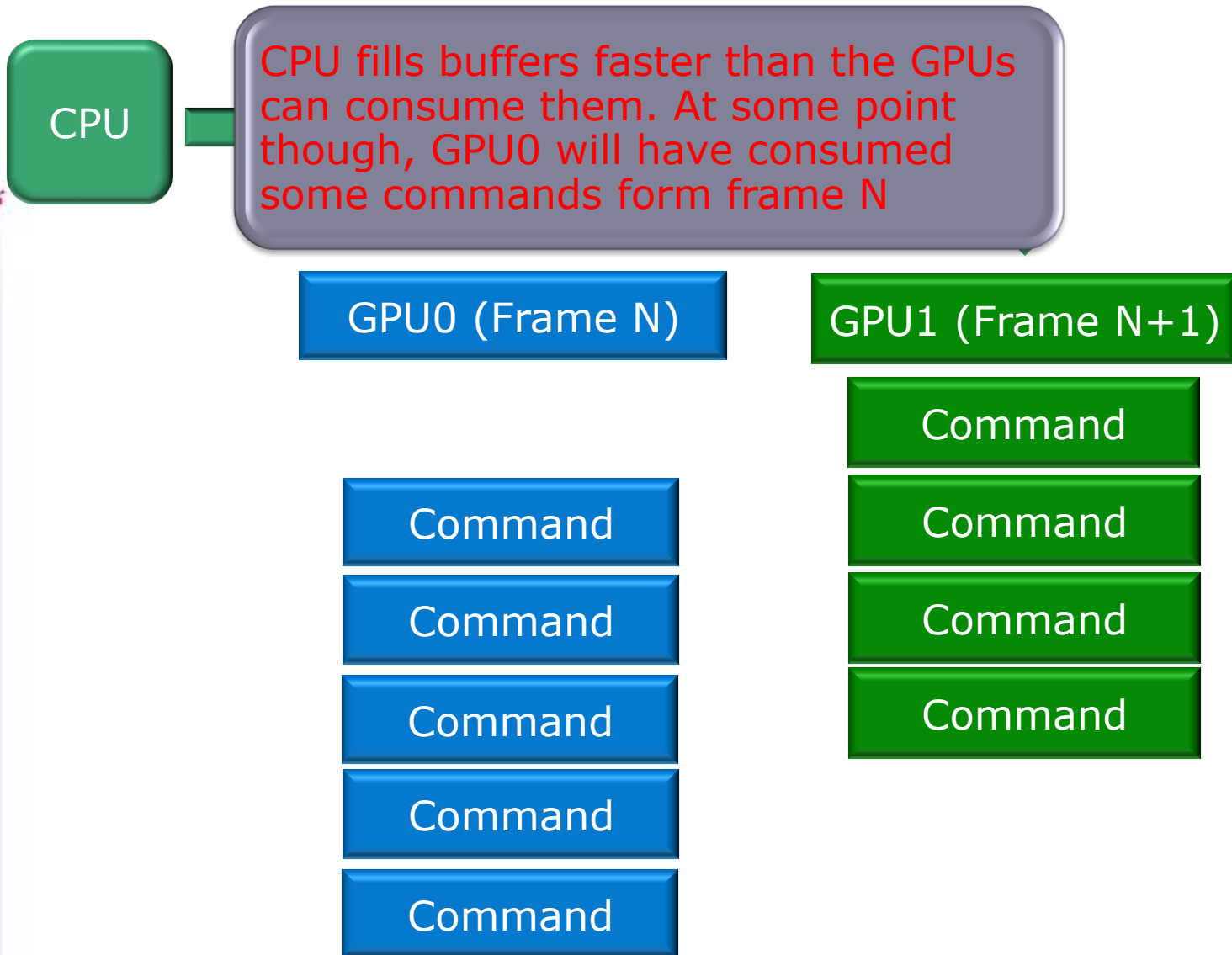


How does AFR Work?



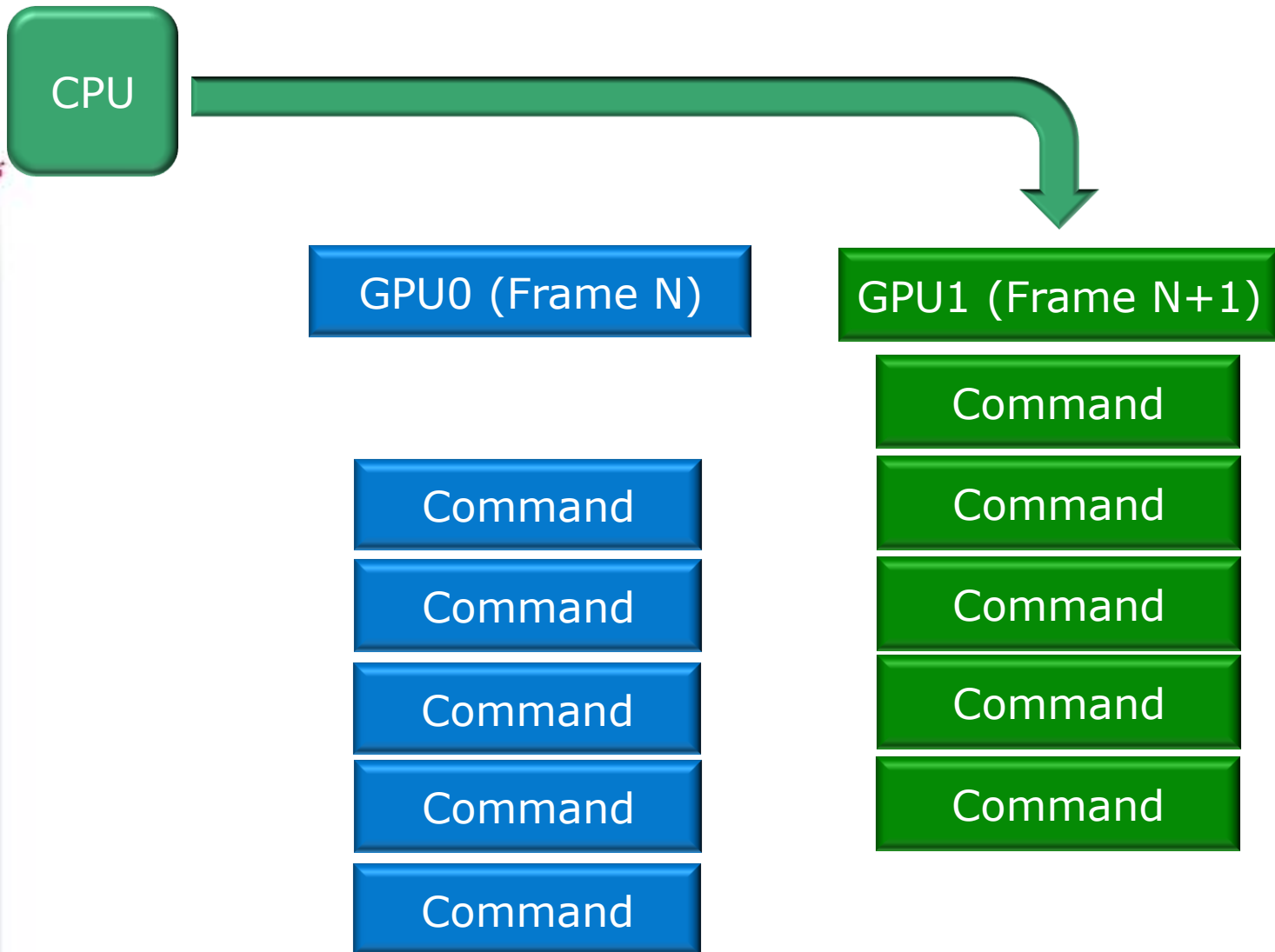


How does AFR Work?



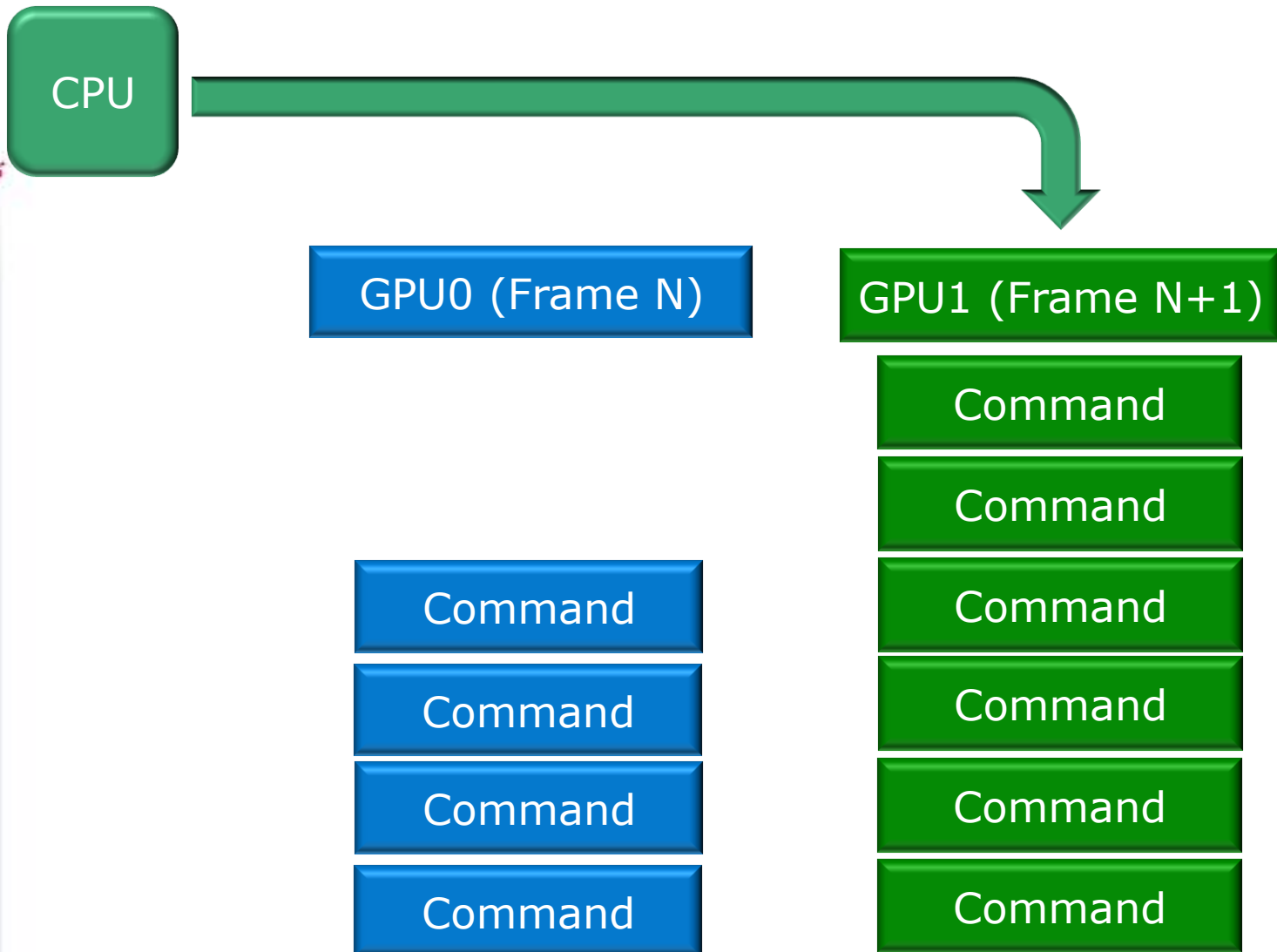


How does AFR Work?





How does AFR Work?





Driver/Programming Considerations



Driver modes for MGPU

④ Compatible AFR Mode

- ④ Default mode
- ④ Driver checks for AFR unfriendly behaviour
- ④ Will work around AFR unfriendly behaviour

④ Full AFR Mode (Application Profile)

- ④ Driver recognises EXE name
- ④ Behaviour fully guided by profile
- ④ Best performance – no checking
- ④ We need developer input here!
- ④ Rename EXE to "AFR-FriendlyD3D.exe"
 - ④ No checking : Speed & compatibility test



Programming for MGPU

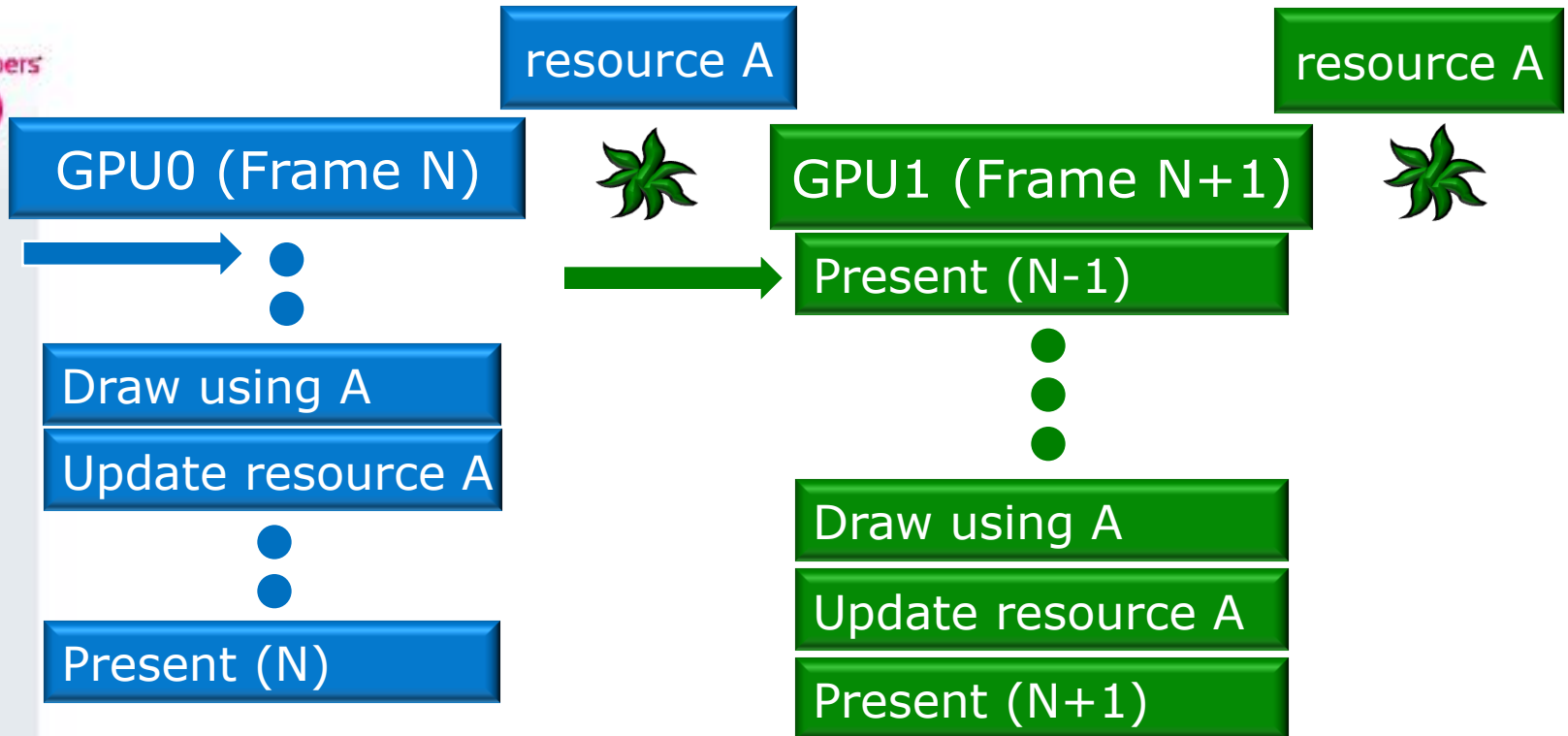
- ⌚ Current MGPU setups are not shared memory architectures
 - ⌚ Apps need to behave well to scale
 - ⌚ Need to keep resources in sync on all GPUs
- ⌚ Know what GPU(s) are rendering the current frame
 - ⌚ Critical to adapt application behavior
- ⌚ Use AMD / NVIDIA libraries
 - ⌚ **AMD: Query the number of GPUs**
 - ⌚ **Nvidia: Query MGPU topology**



Common Pitfalls & Solutions

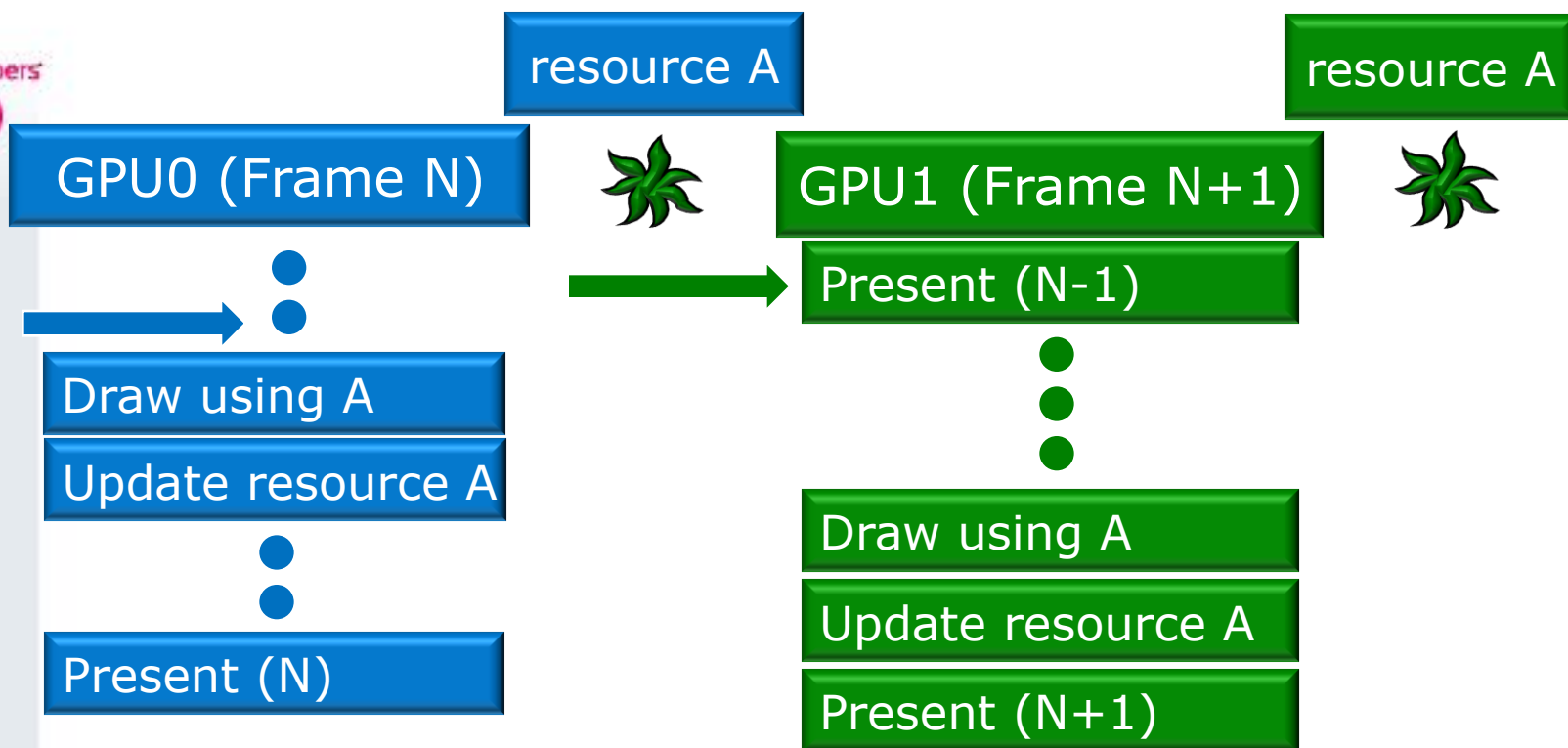


Pitfall: Dependencies Between Frames



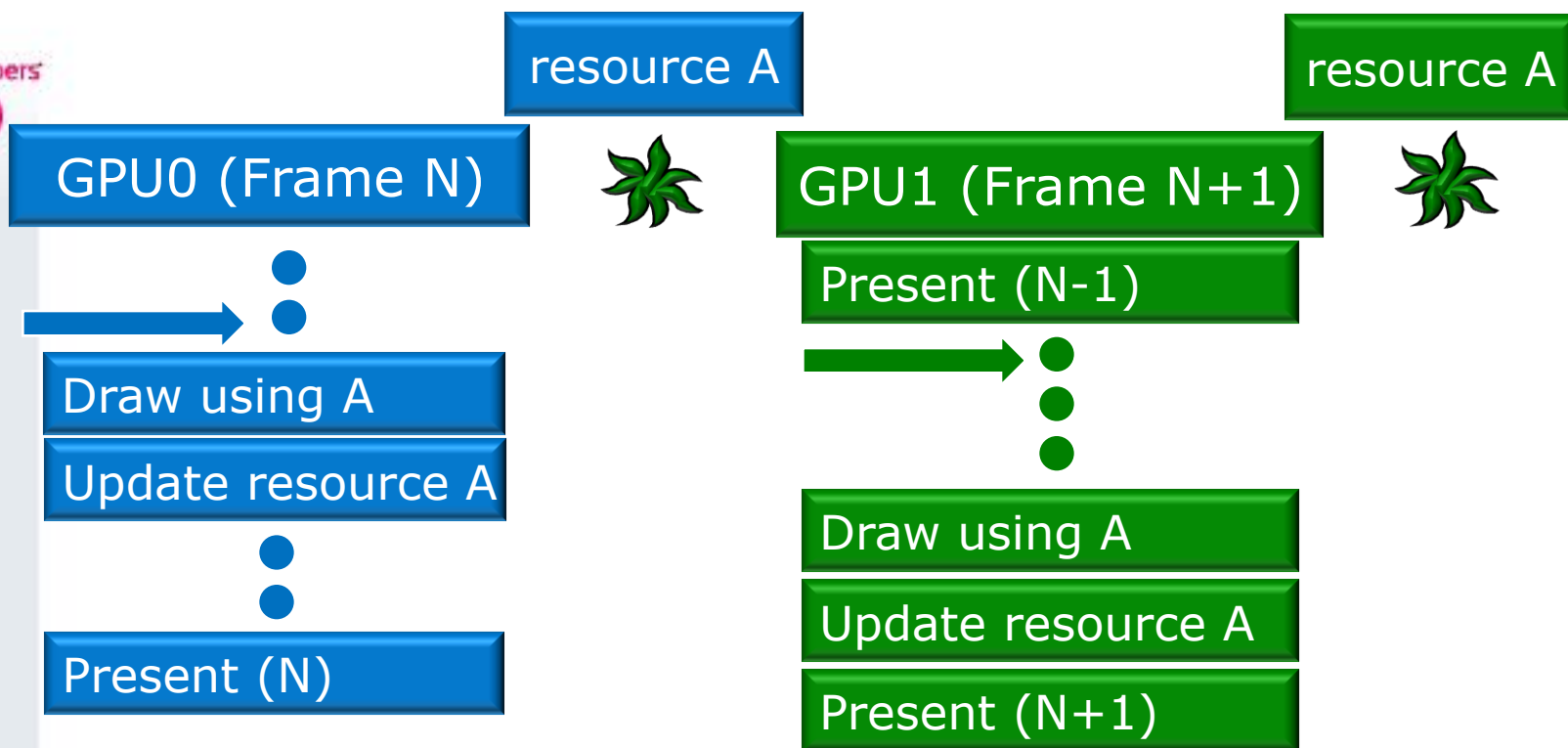


Pitfall: Dependencies Between Frames





Pitfall: Dependencies Between Frames

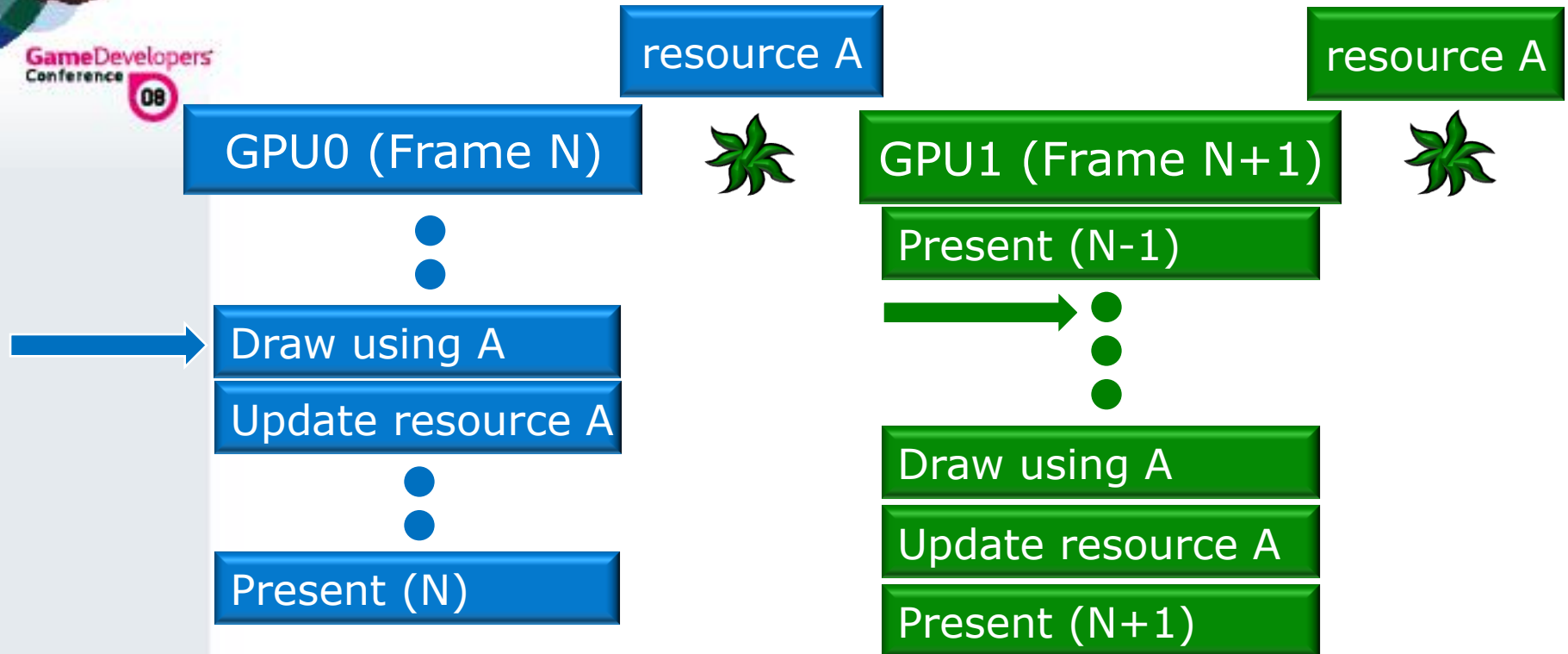




Game Developers
Conference

08

Pitfall: Dependencies Between Frames

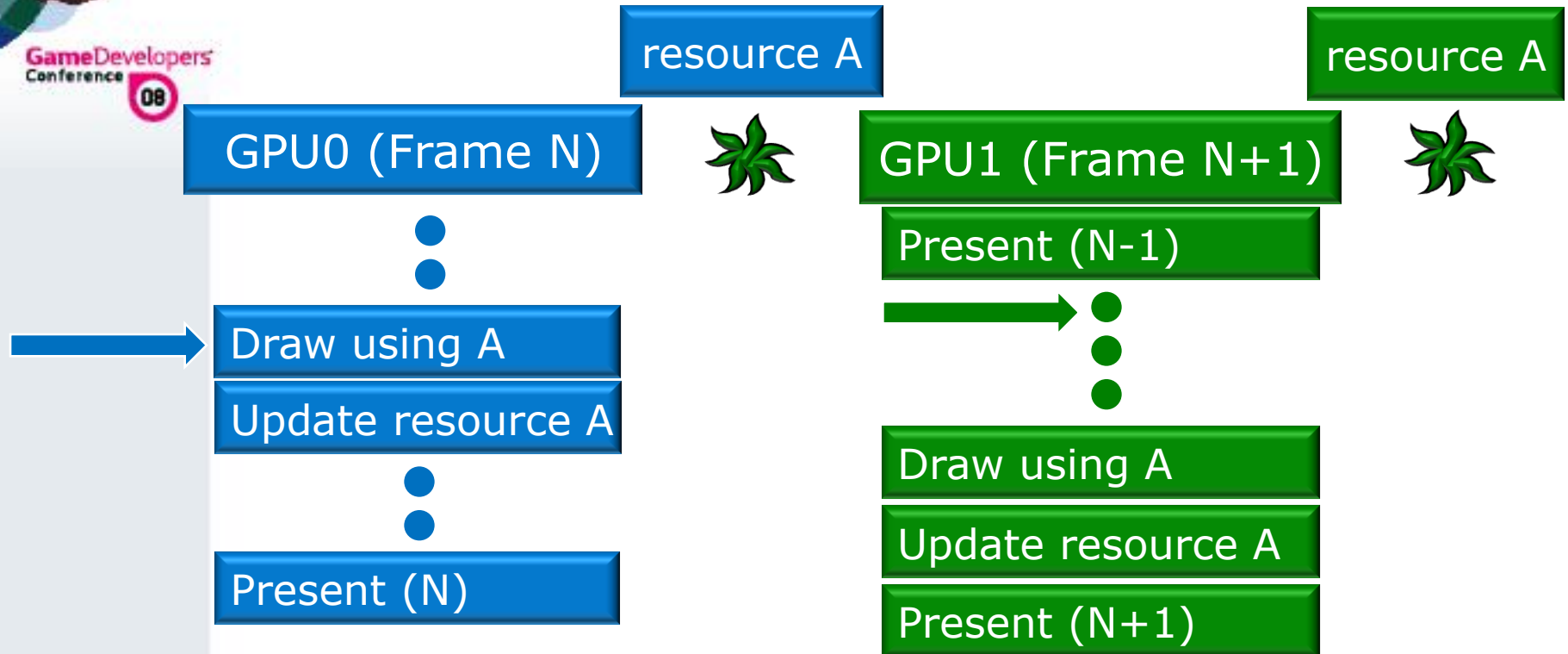




Game Developers
Conference

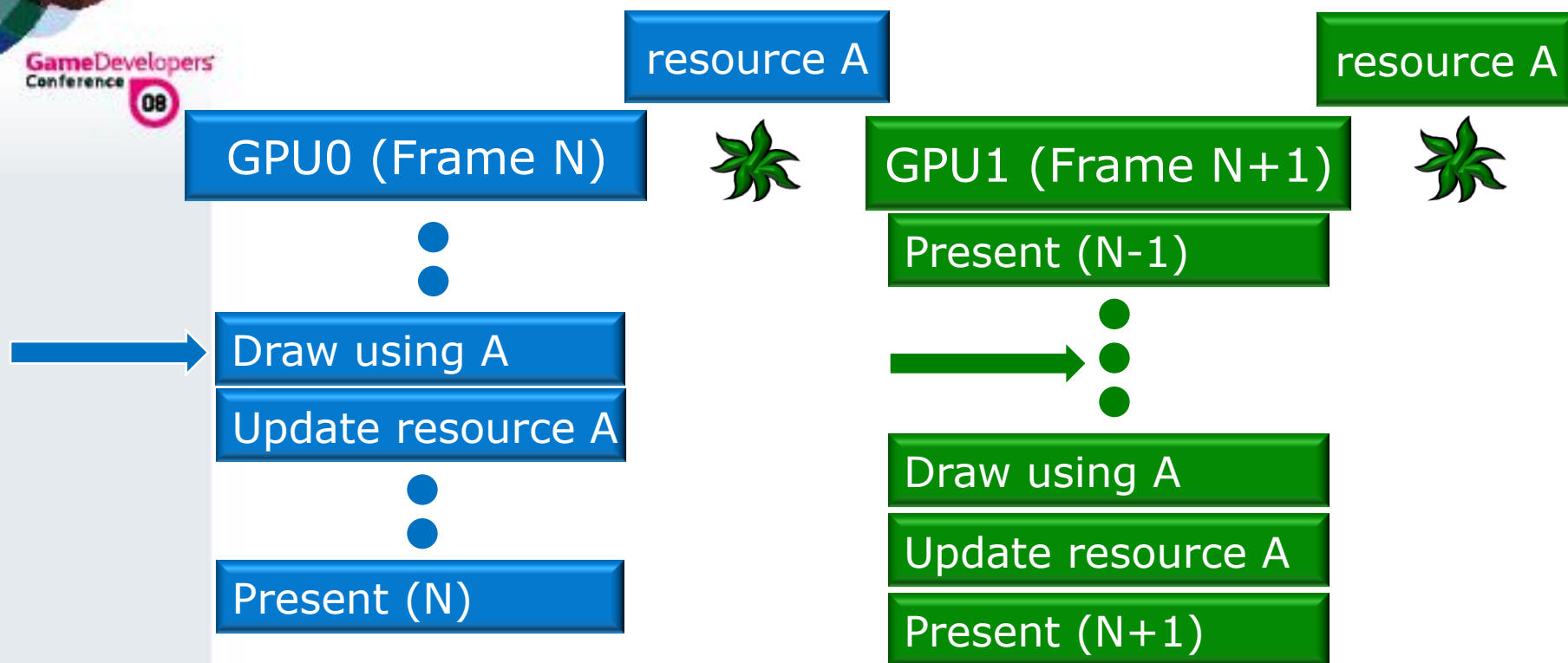
08

Pitfall: Dependencies Between Frames

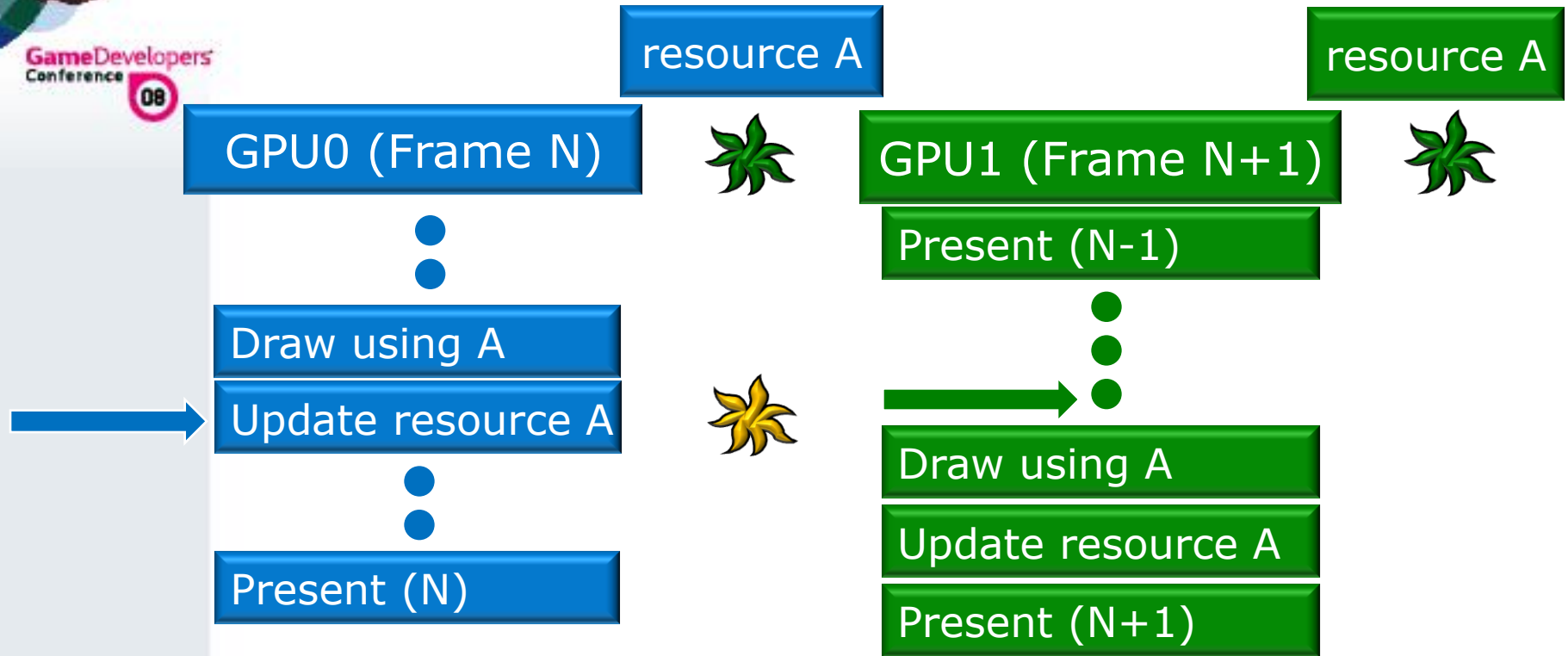




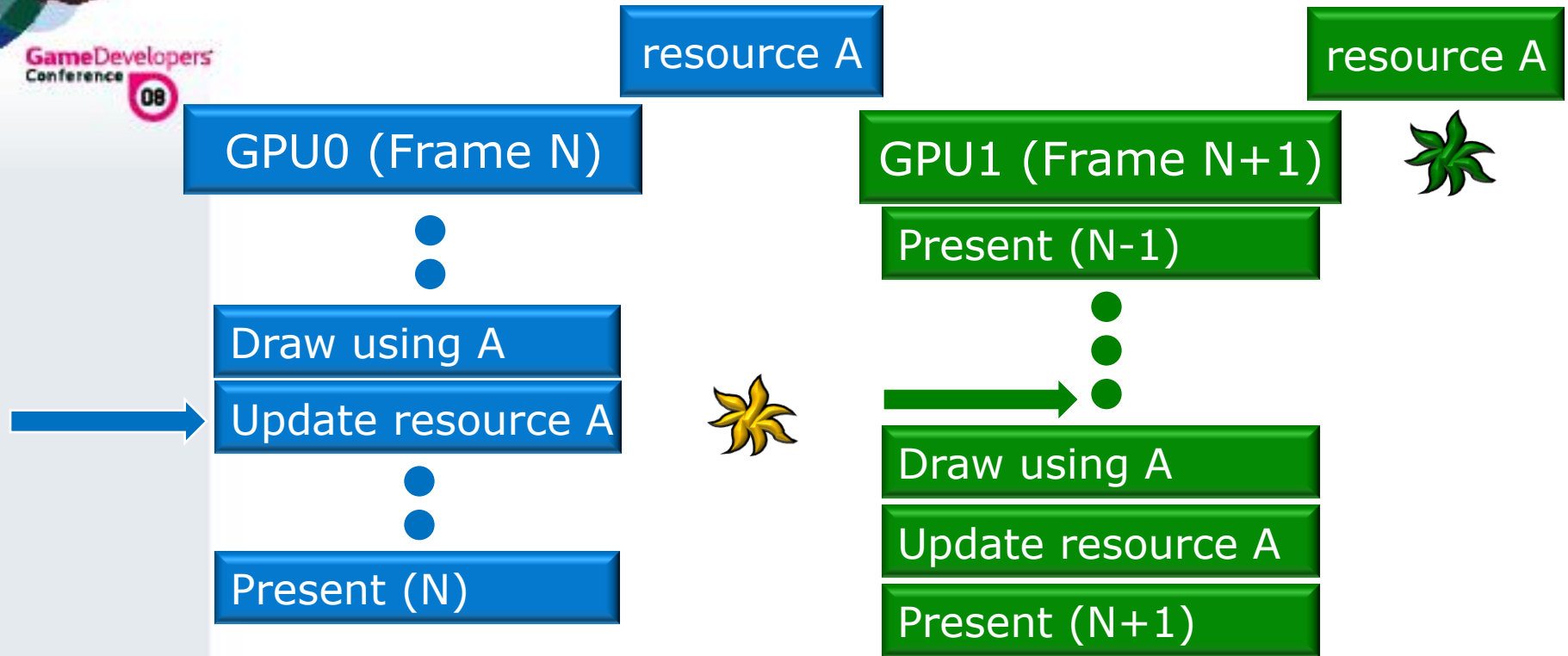
Pitfall: Dependencies Between Frames



Pitfall: Dependencies Between Frames



Pitfall: Dependencies Between Frames

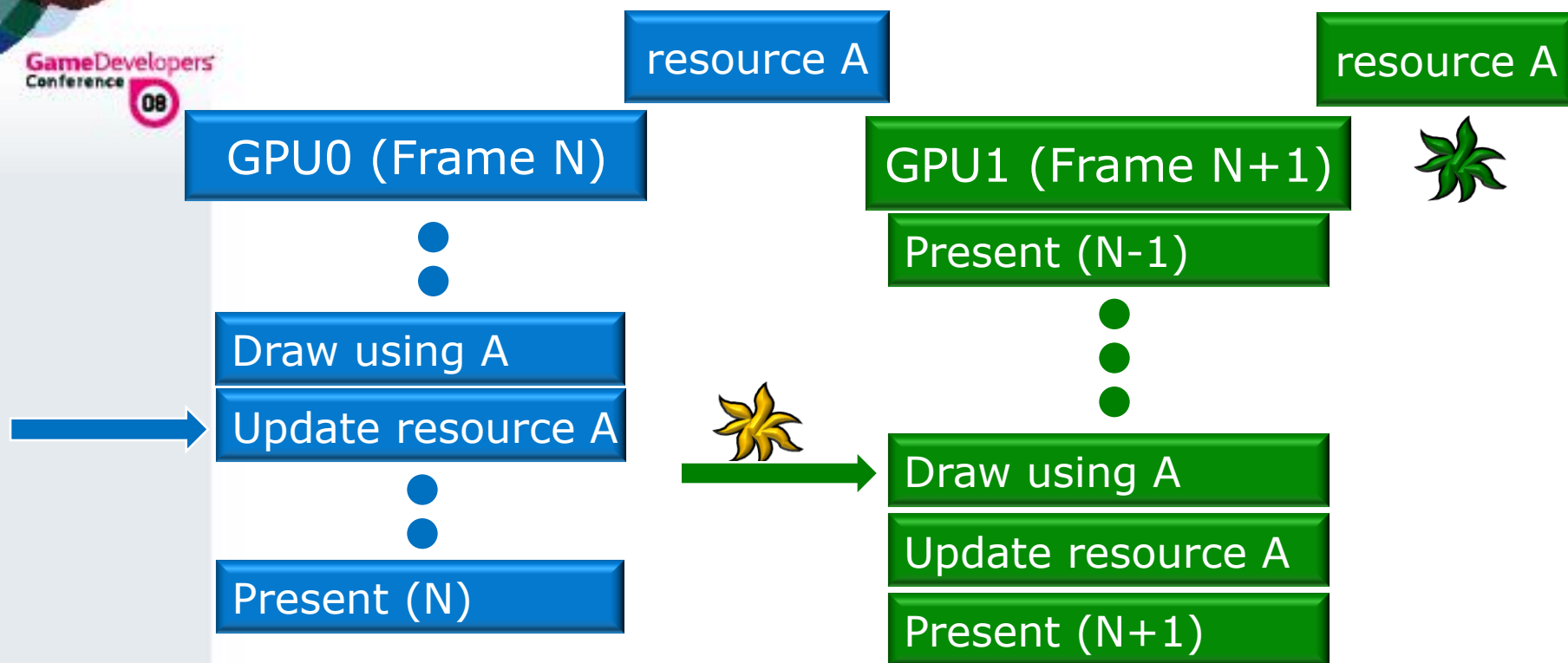




Game Developers
Conference

08

Pitfall: Dependencies Between Frames

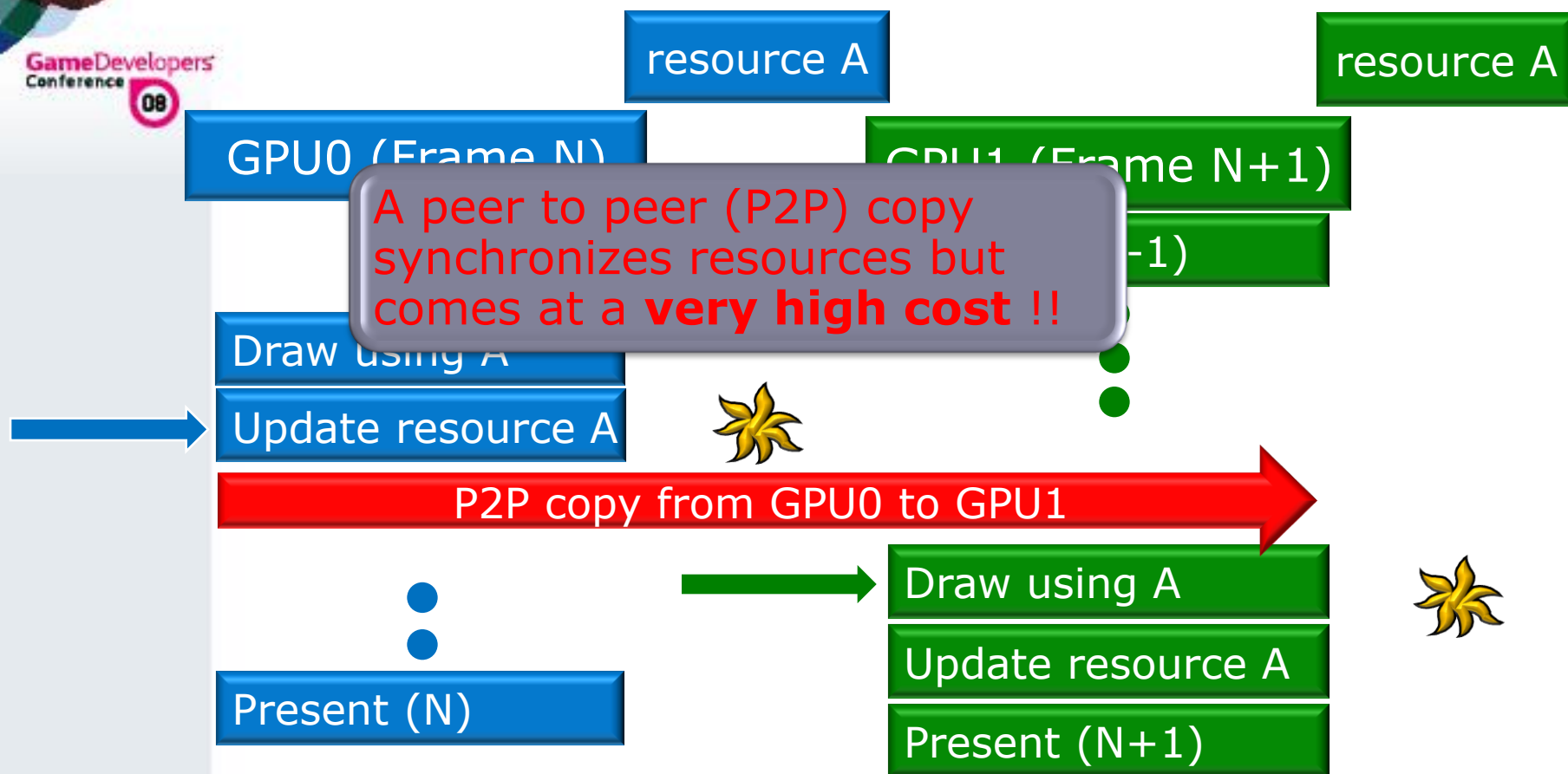


CMP

United Business Media



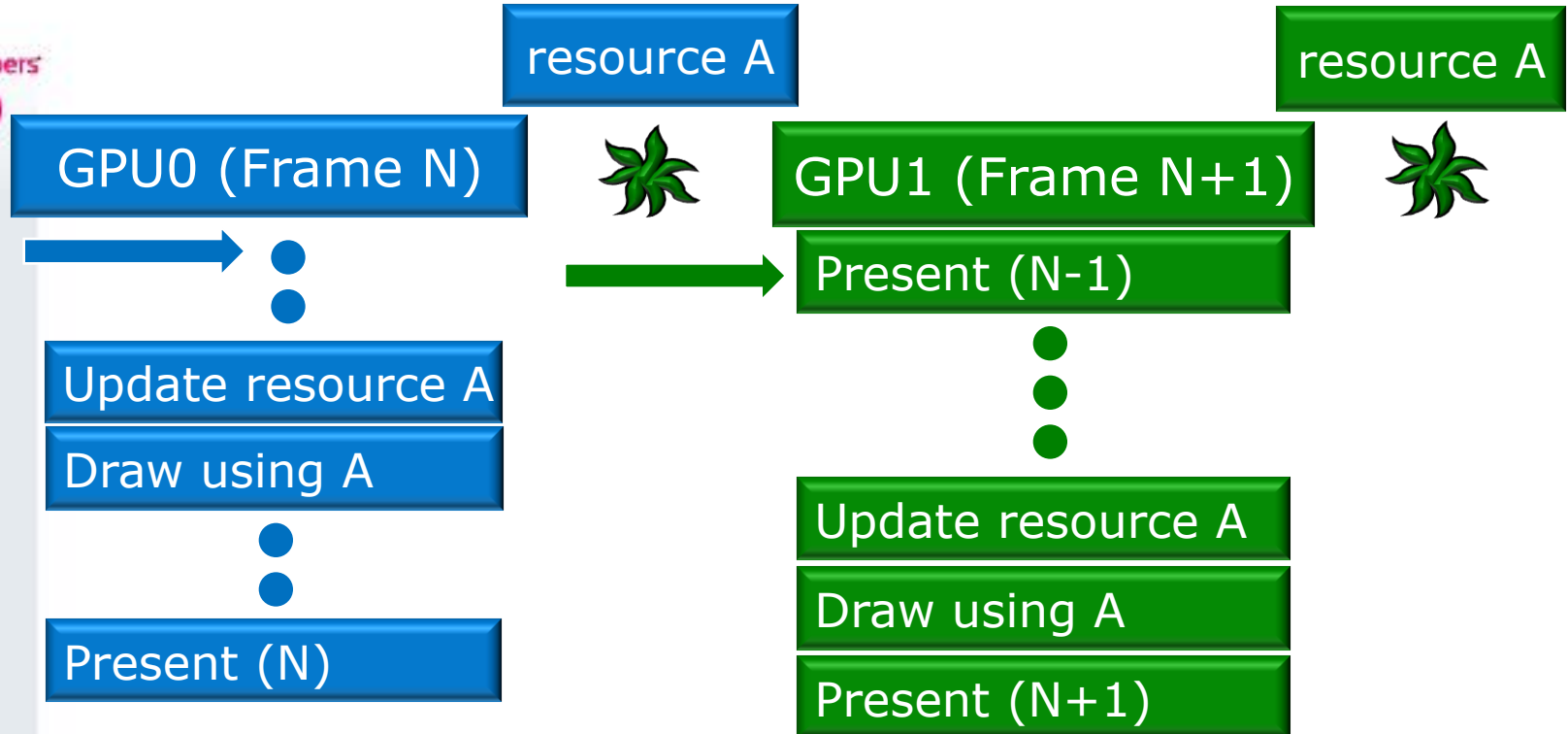
Pitfall: Dependencies Between Frames



Solution: Resources that Change Every Frame

Game Developers
Conference

08



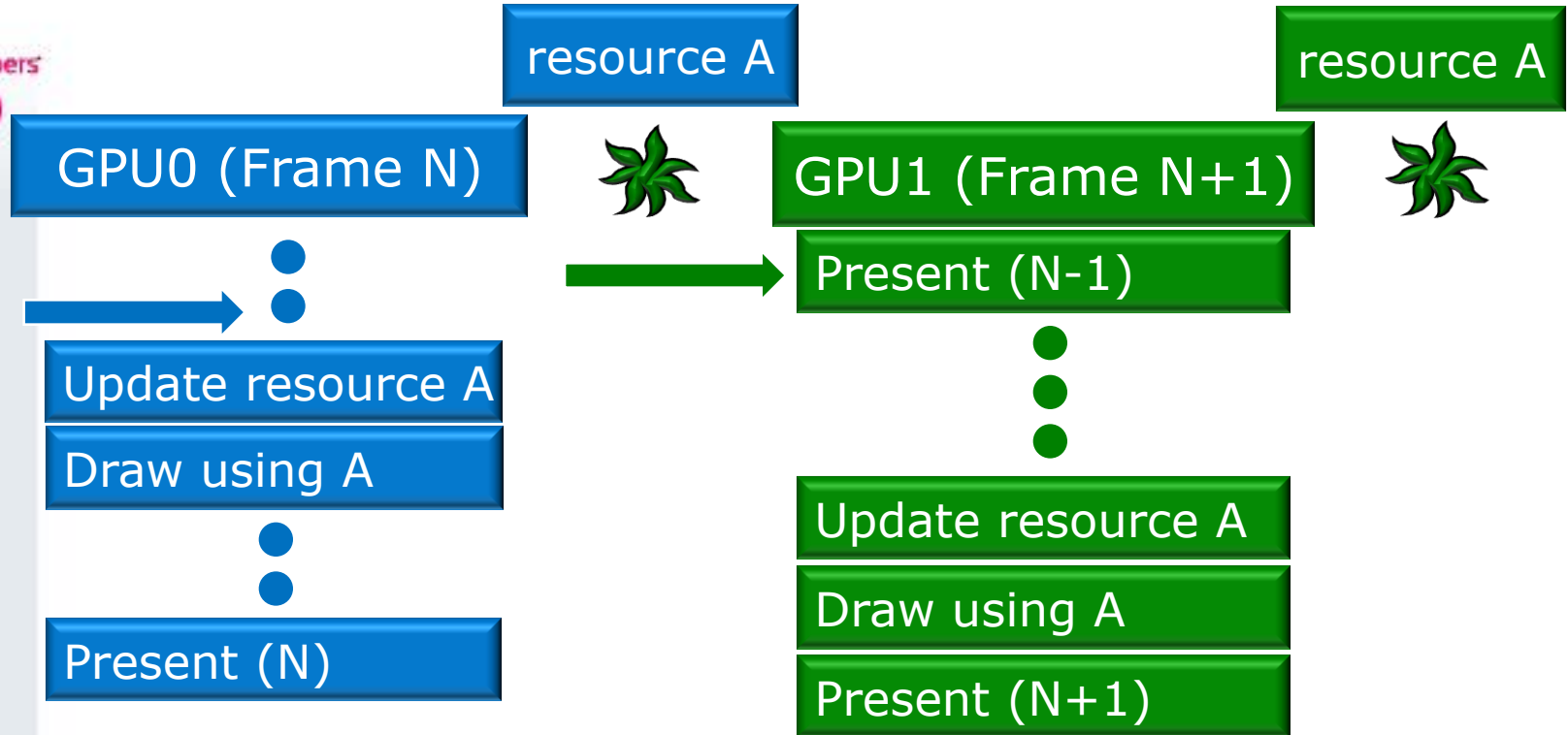
CMP

United Business Media

Solution: Resources that Change Every Frame

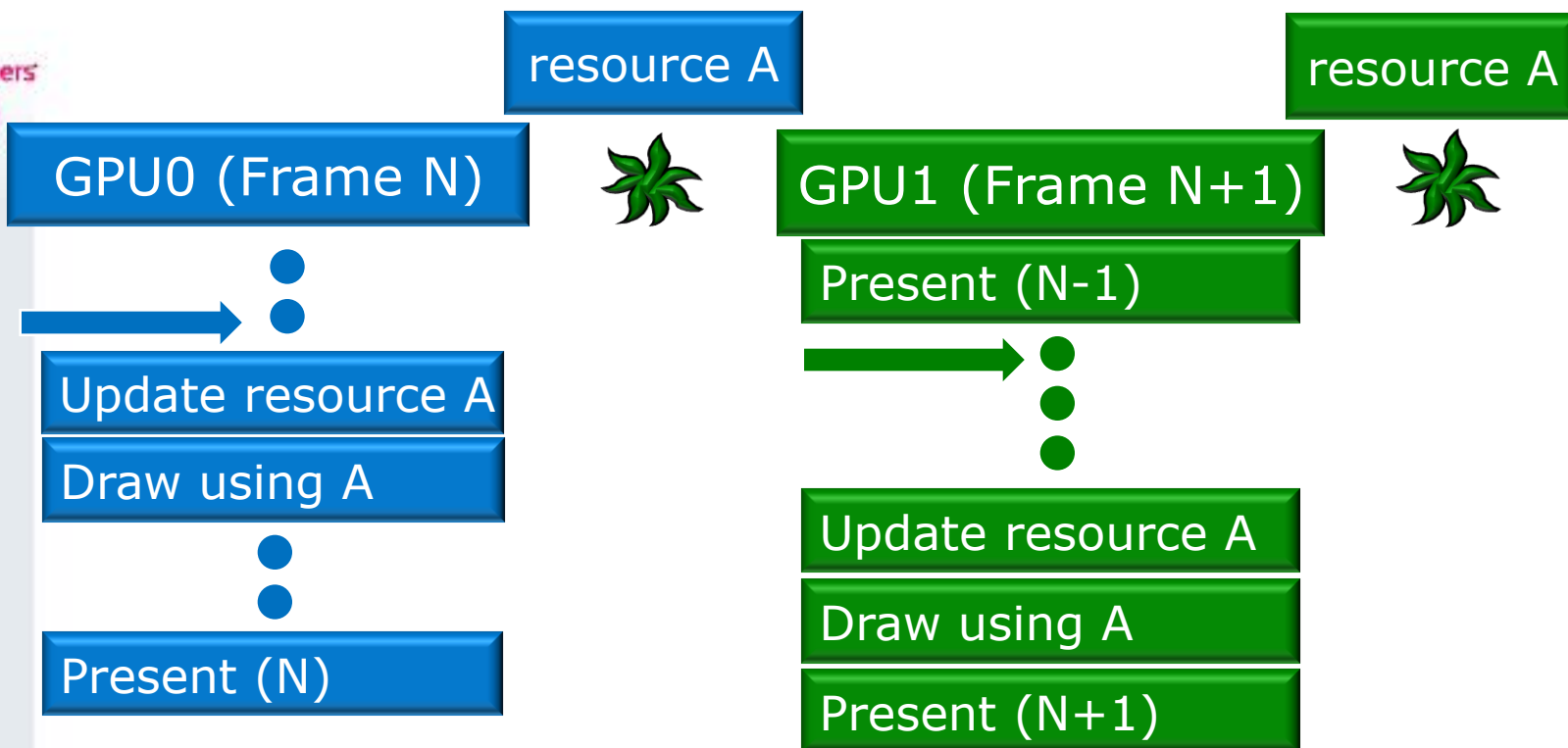
Game Developers
Conference

08



CMP

United Business Media

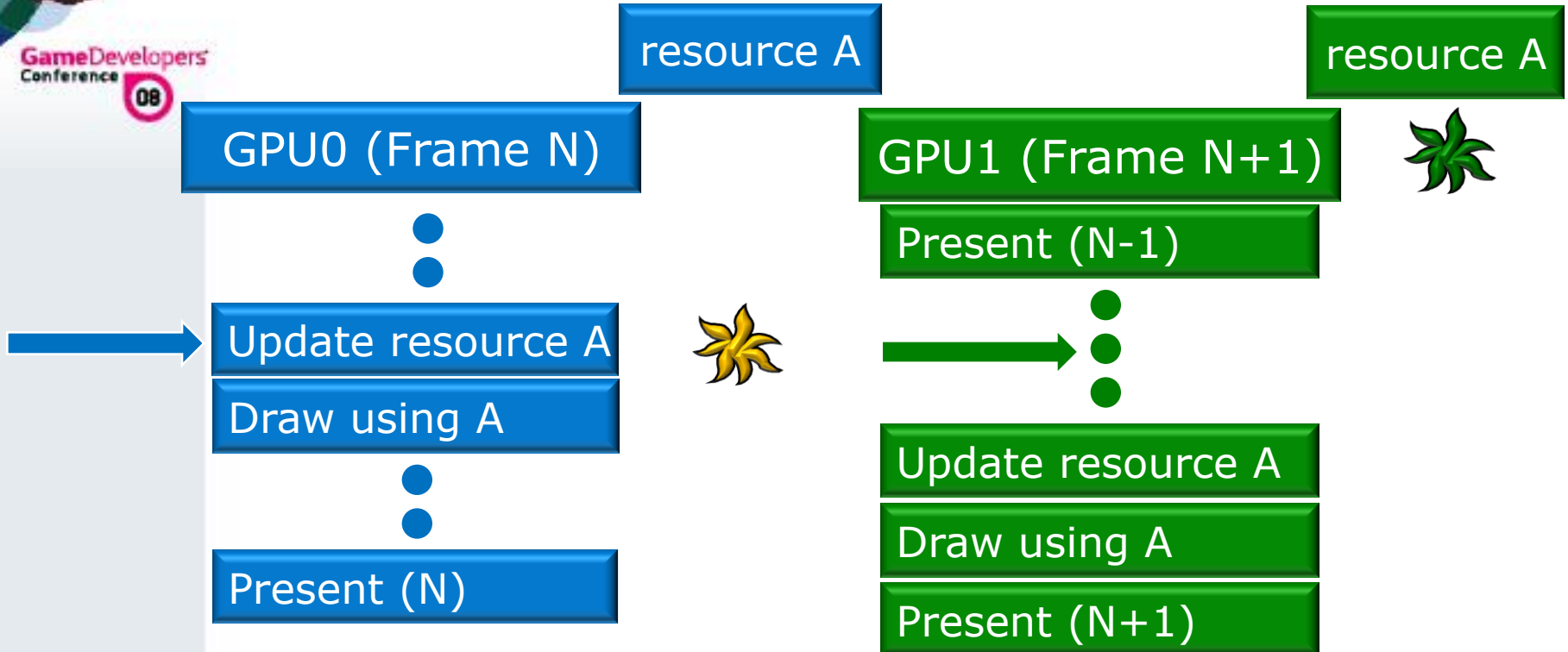




Game Developers
Conference

08

Solution: Resources that Change Every Frame

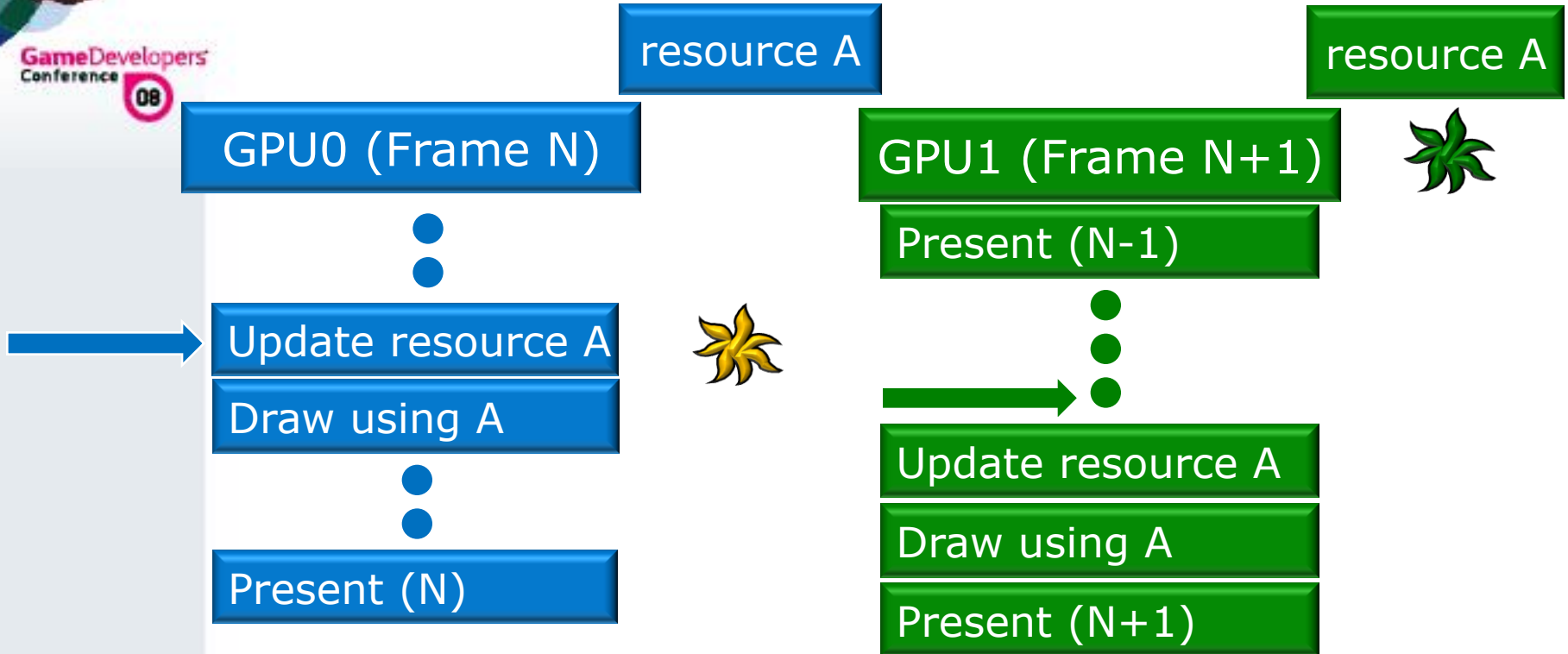




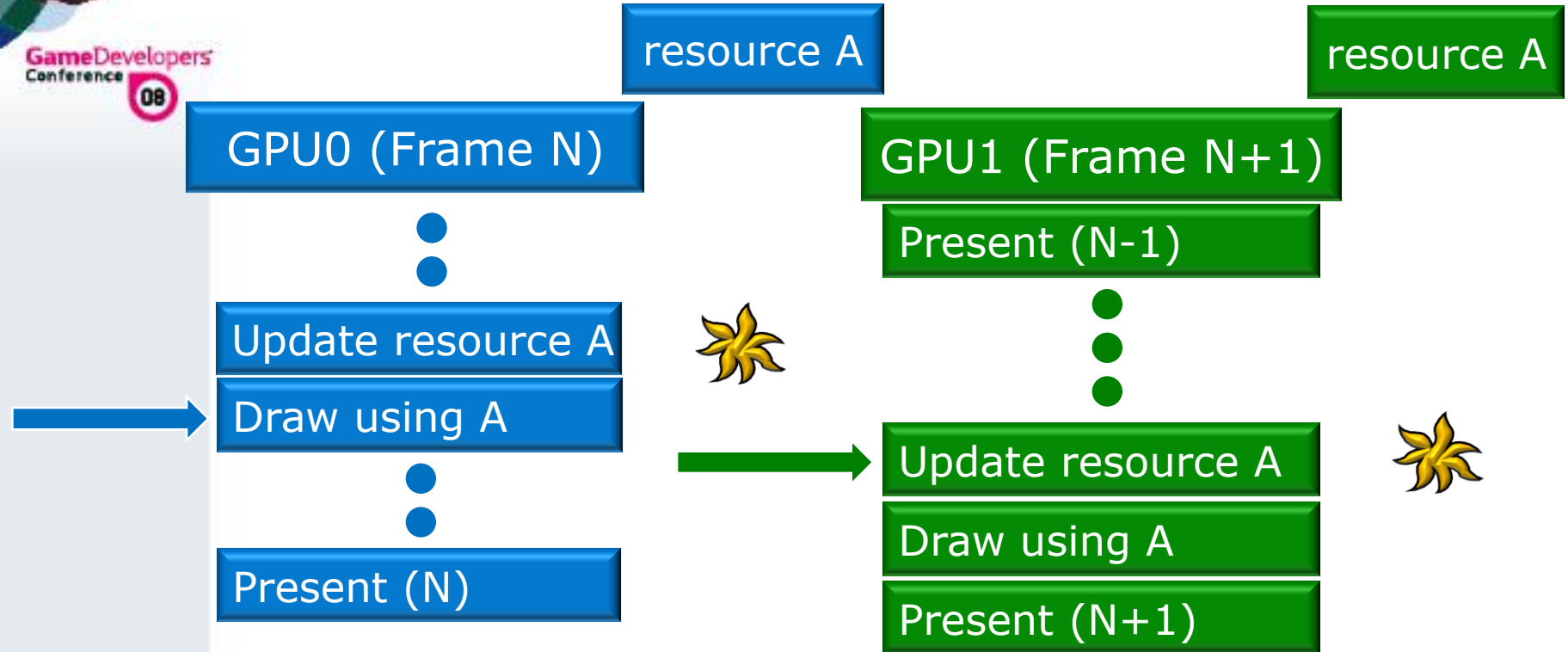
Game Developers
Conference

08

Solution: Resources that Change Every Frame



Solution: Resources that Change Every Frame

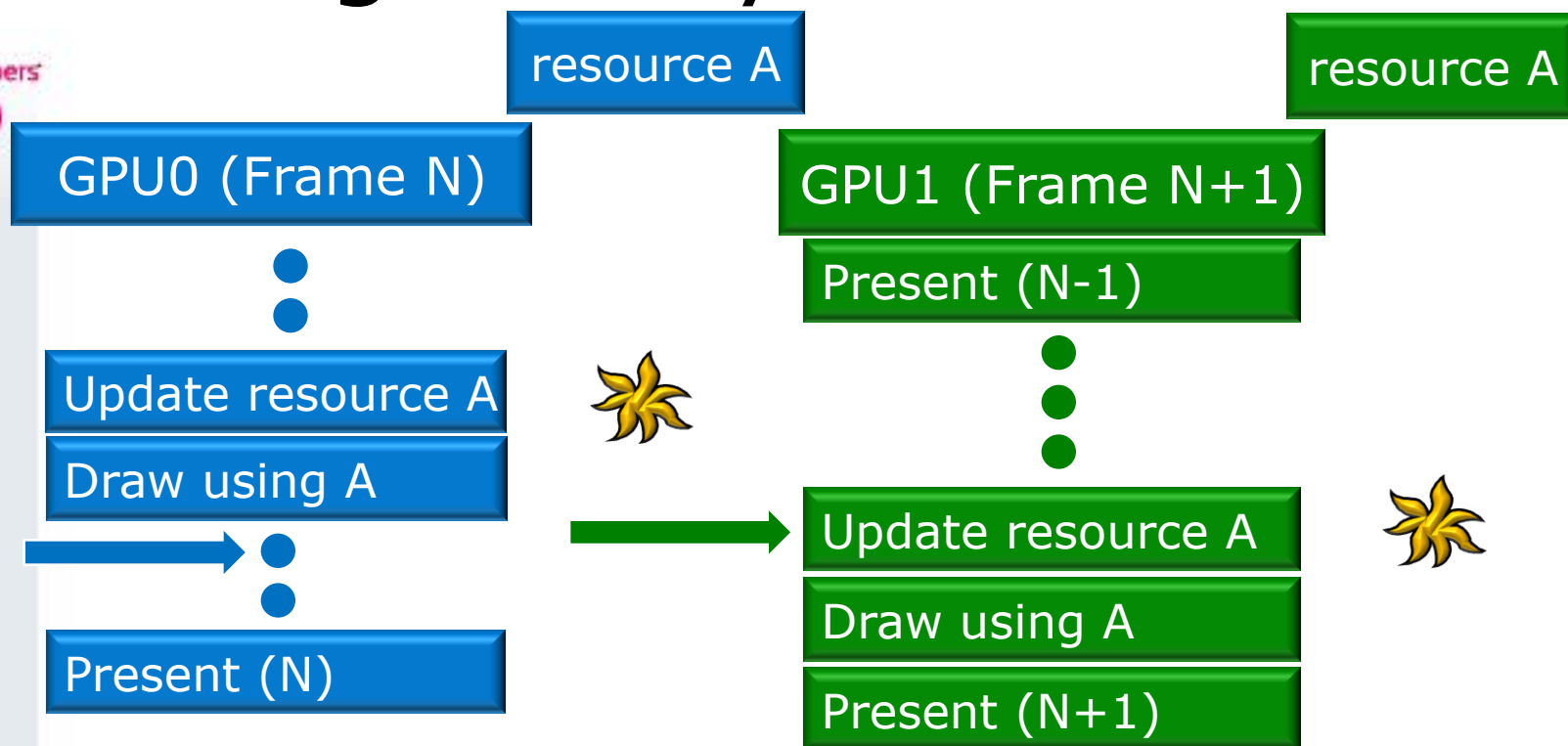




Game Developers
Conference

08

Solution: Resources that Change Every Frame

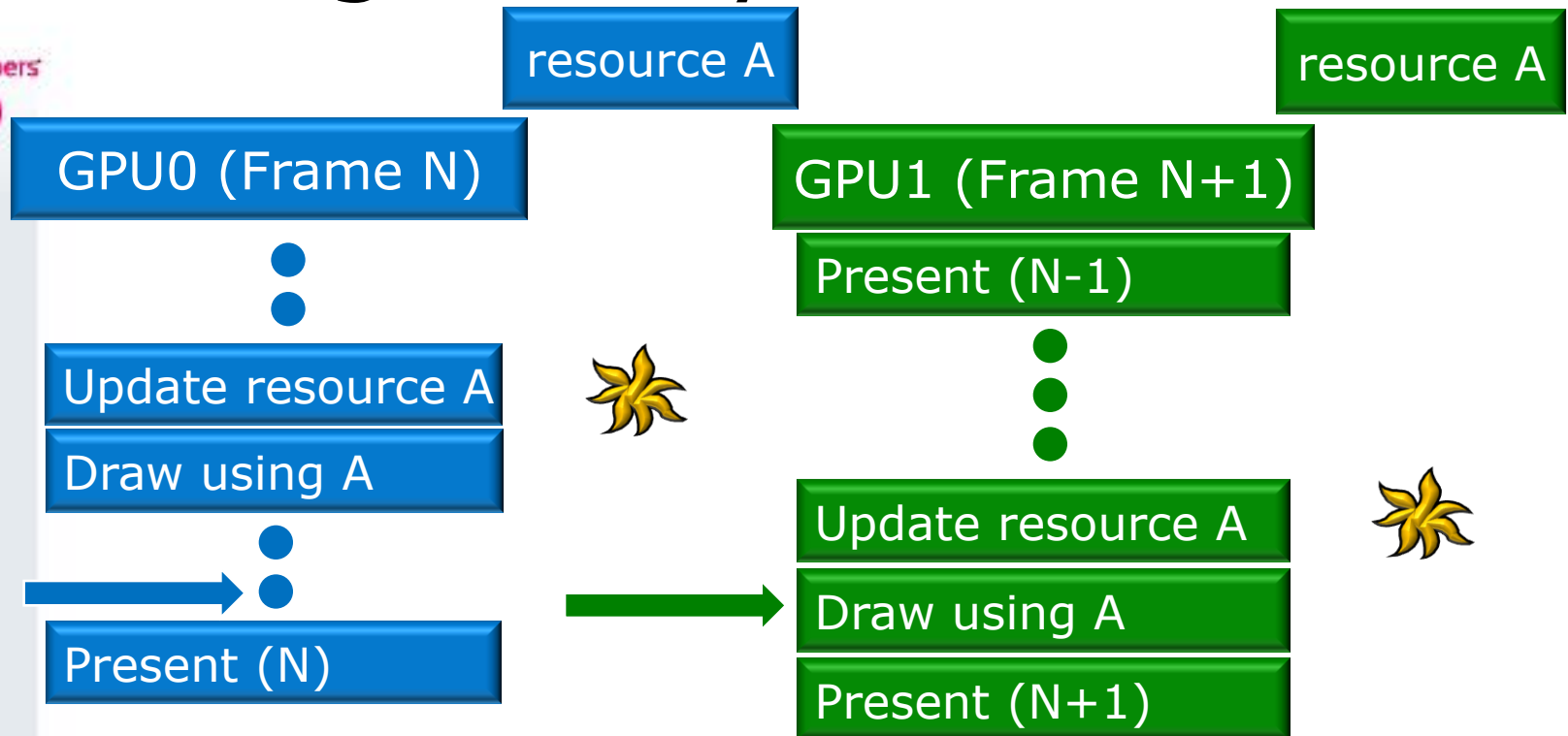




Game Developers
Conference

08

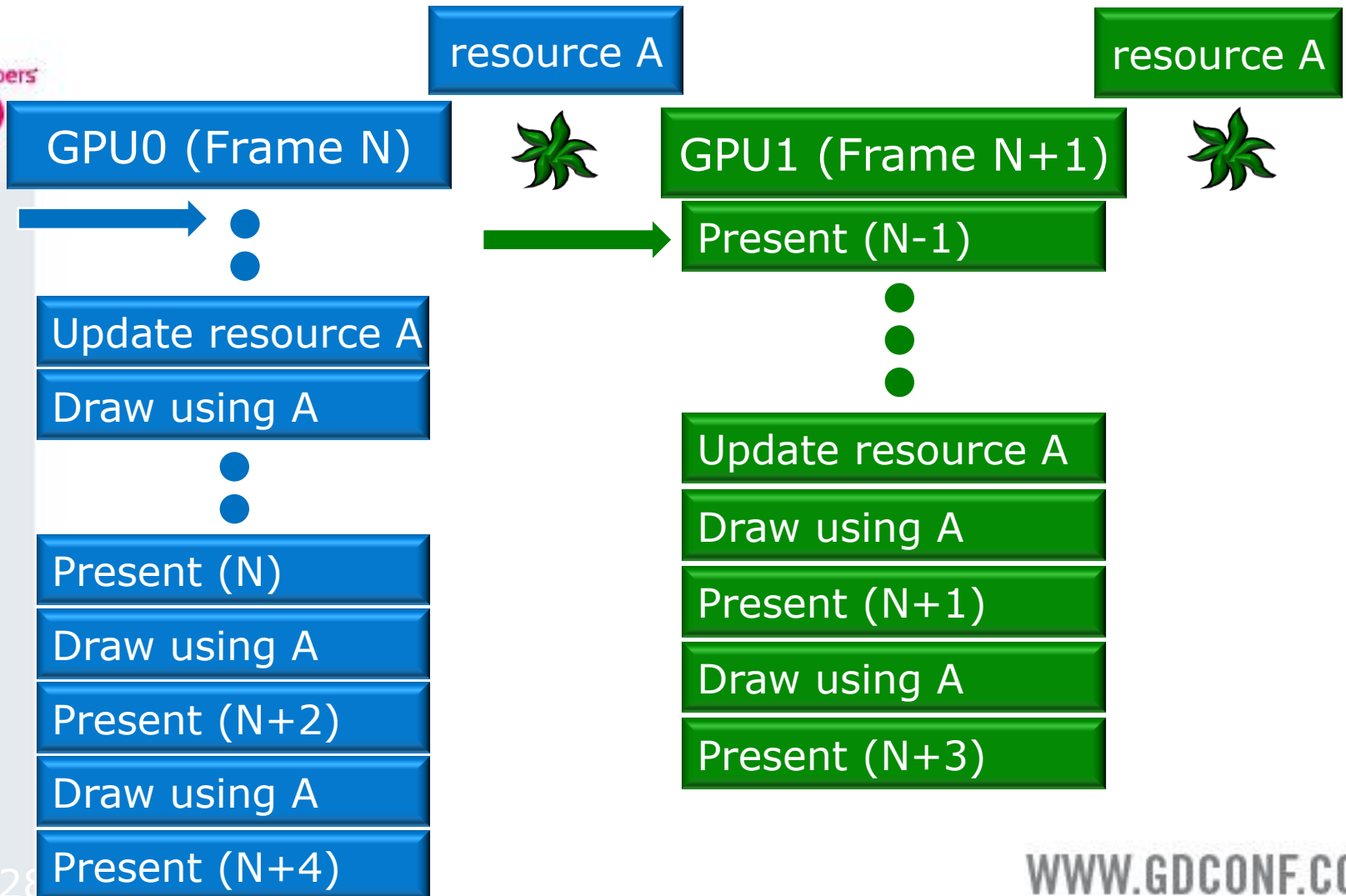
Solution: Resources that Change Every Frame



There are no P2P copies if one always modifies the resource **before** using it within a frame !

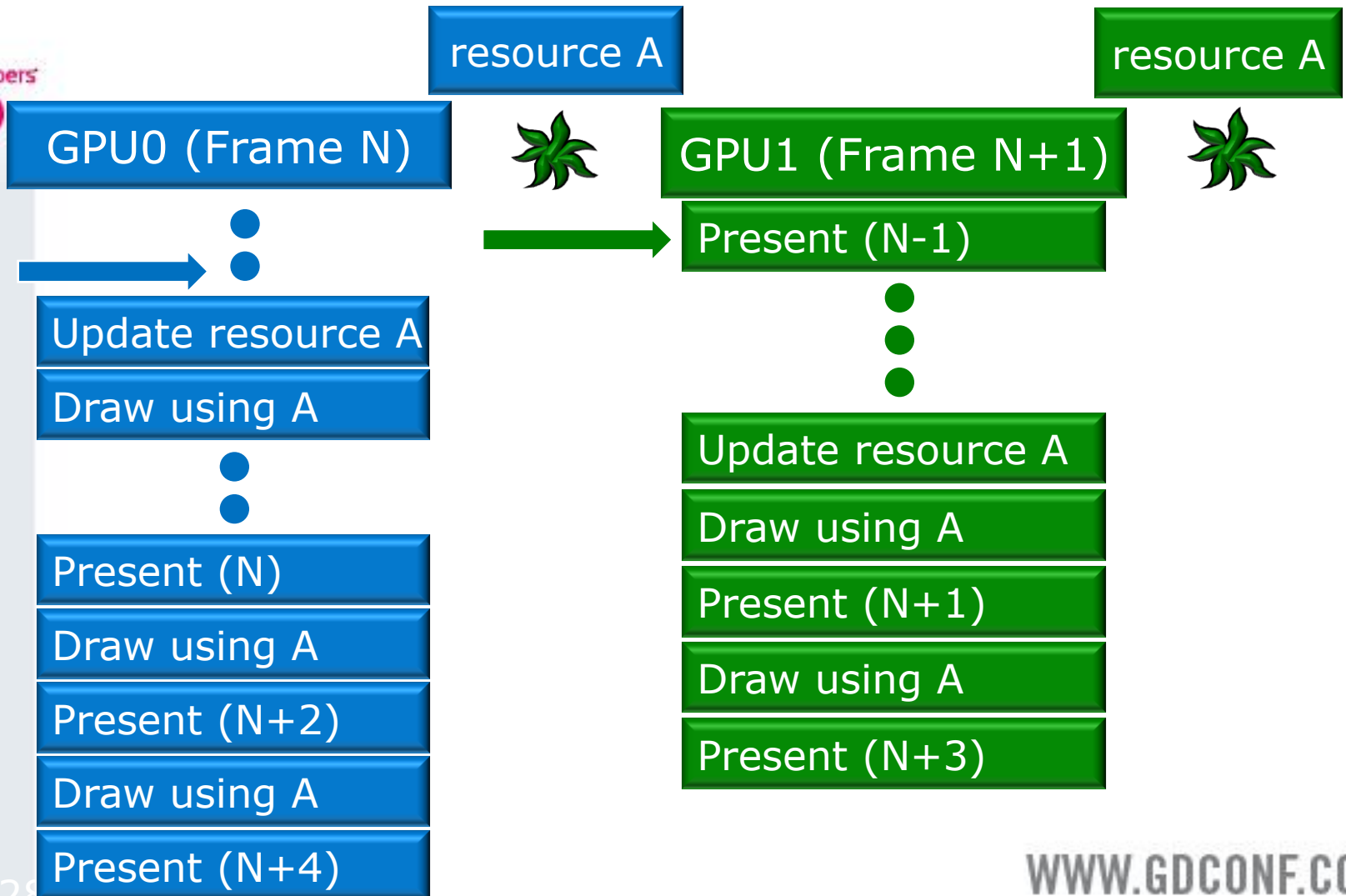


Solution: Resources that Change Every Few Frames





Solution: Resources that Change Every Few Frames

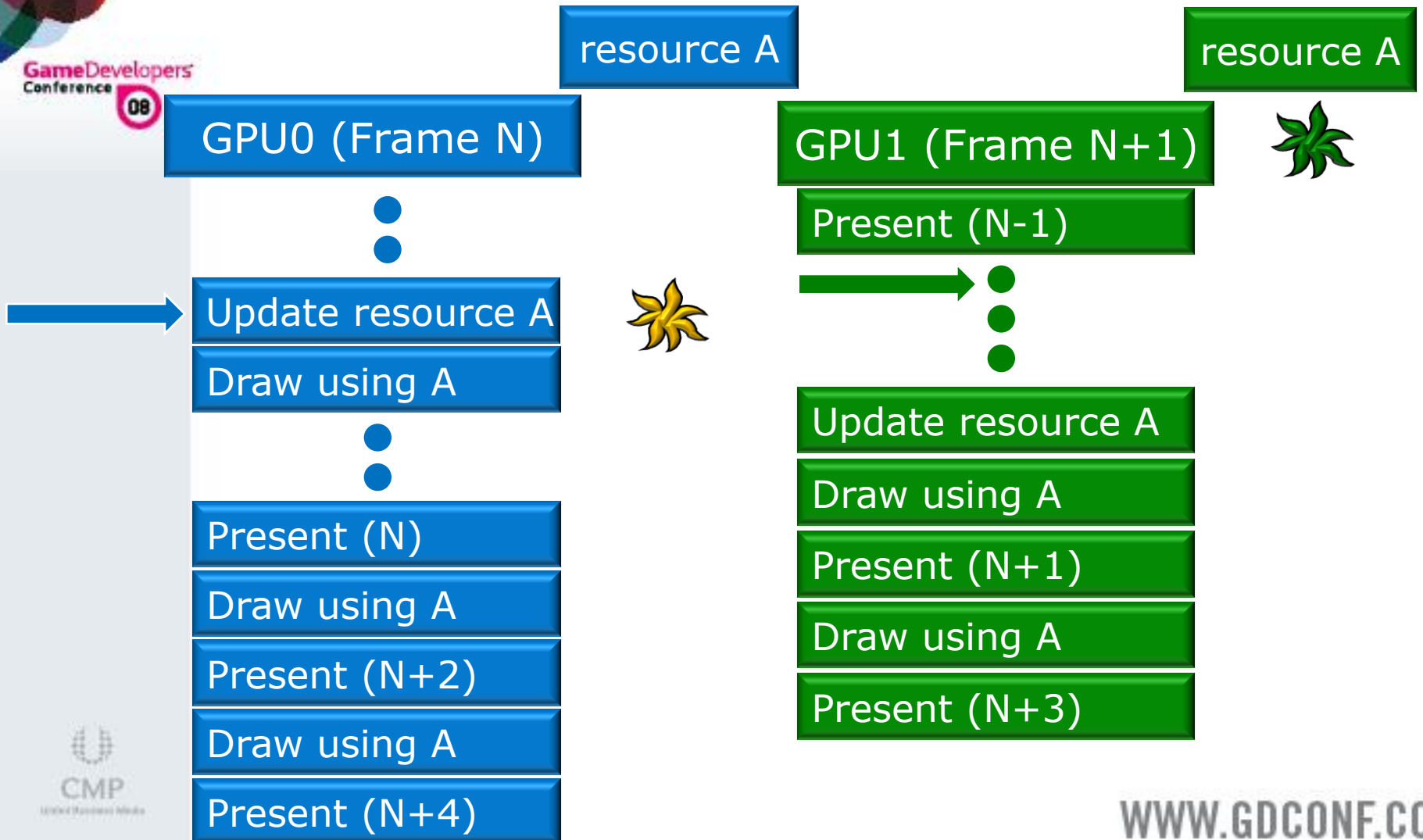




Game Developers
Conference

08

Solution: Resources that Change Every Few Frames

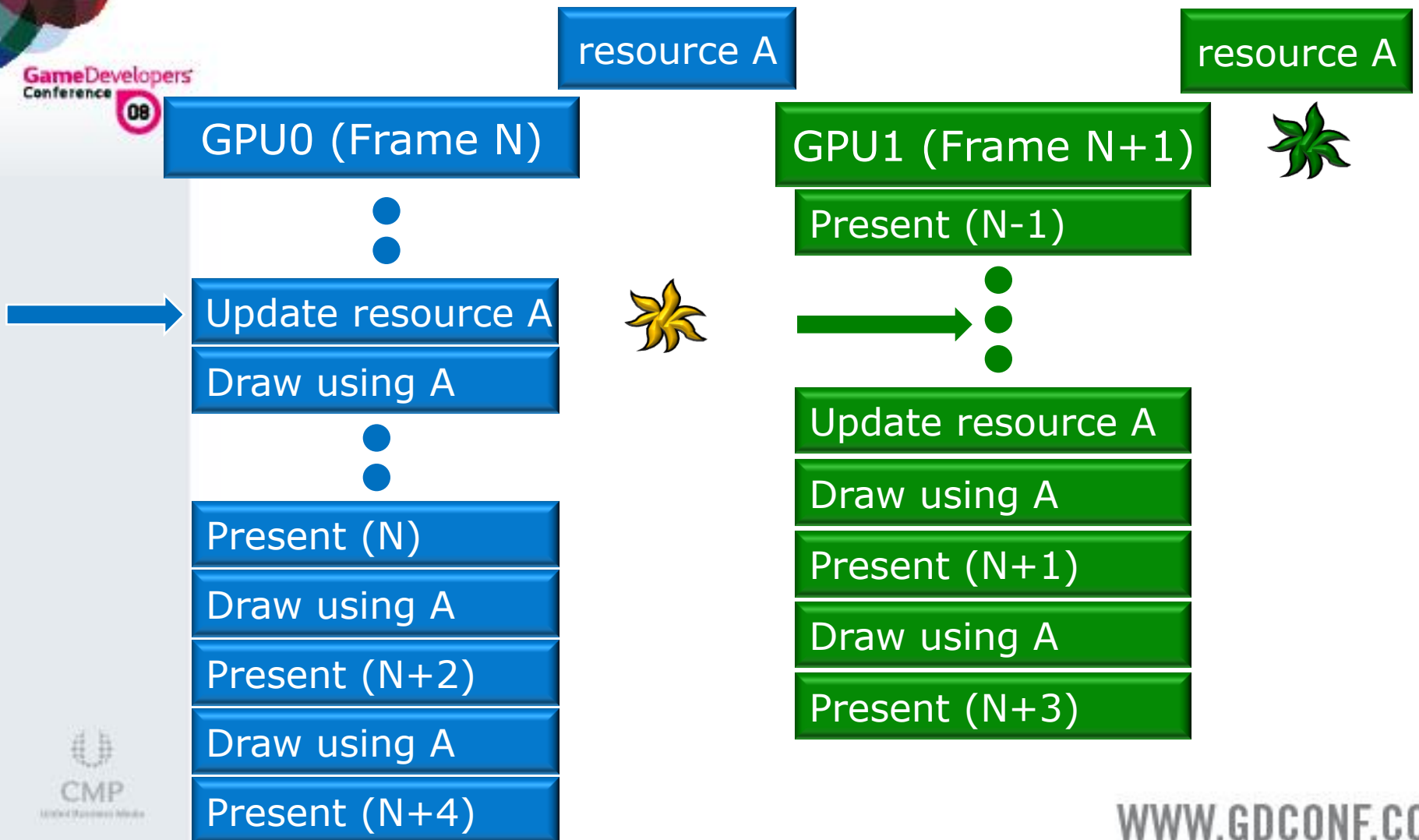




Game Developers
Conference

08

Solution: Resources that Change Every Few Frames

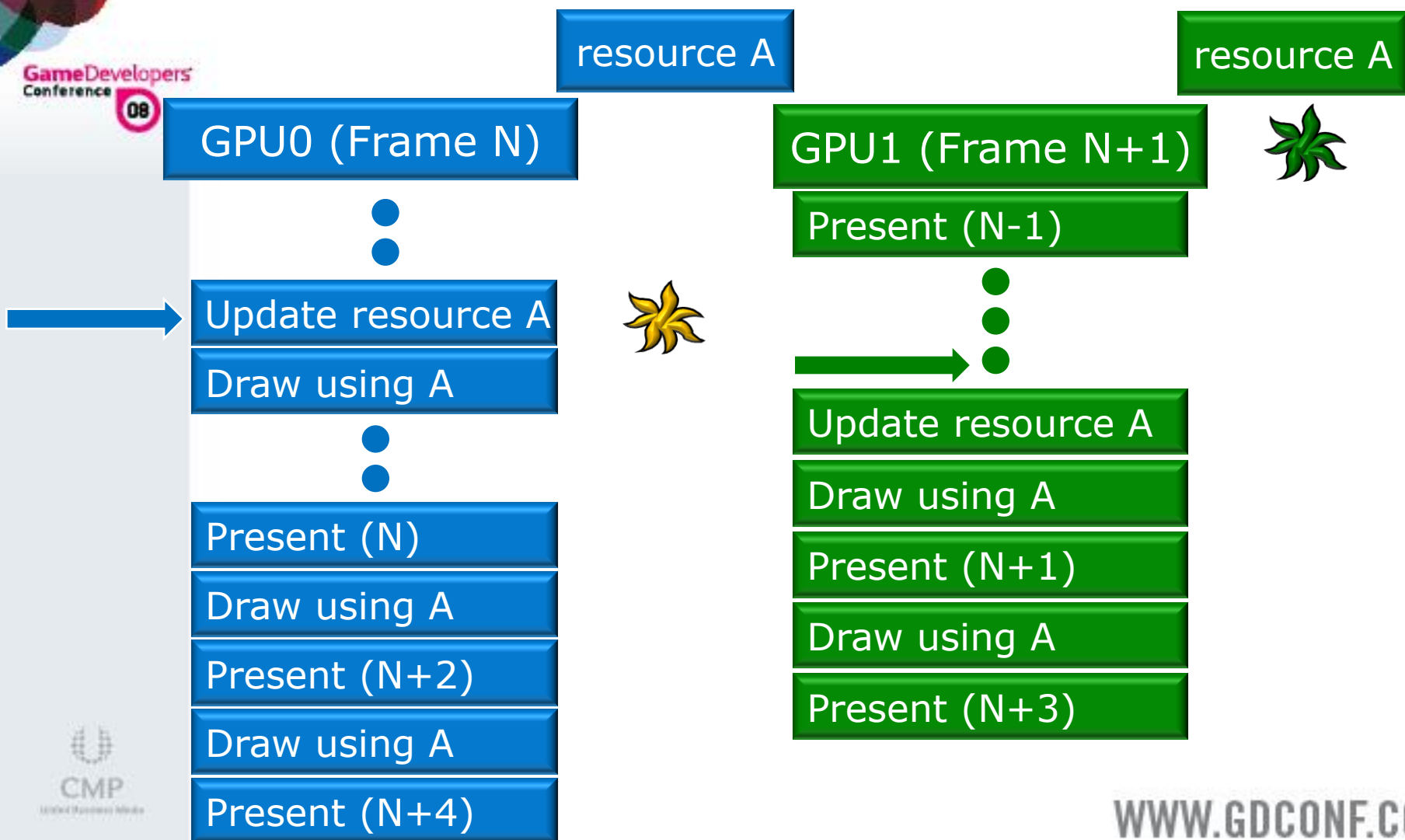




Game Developers
Conference

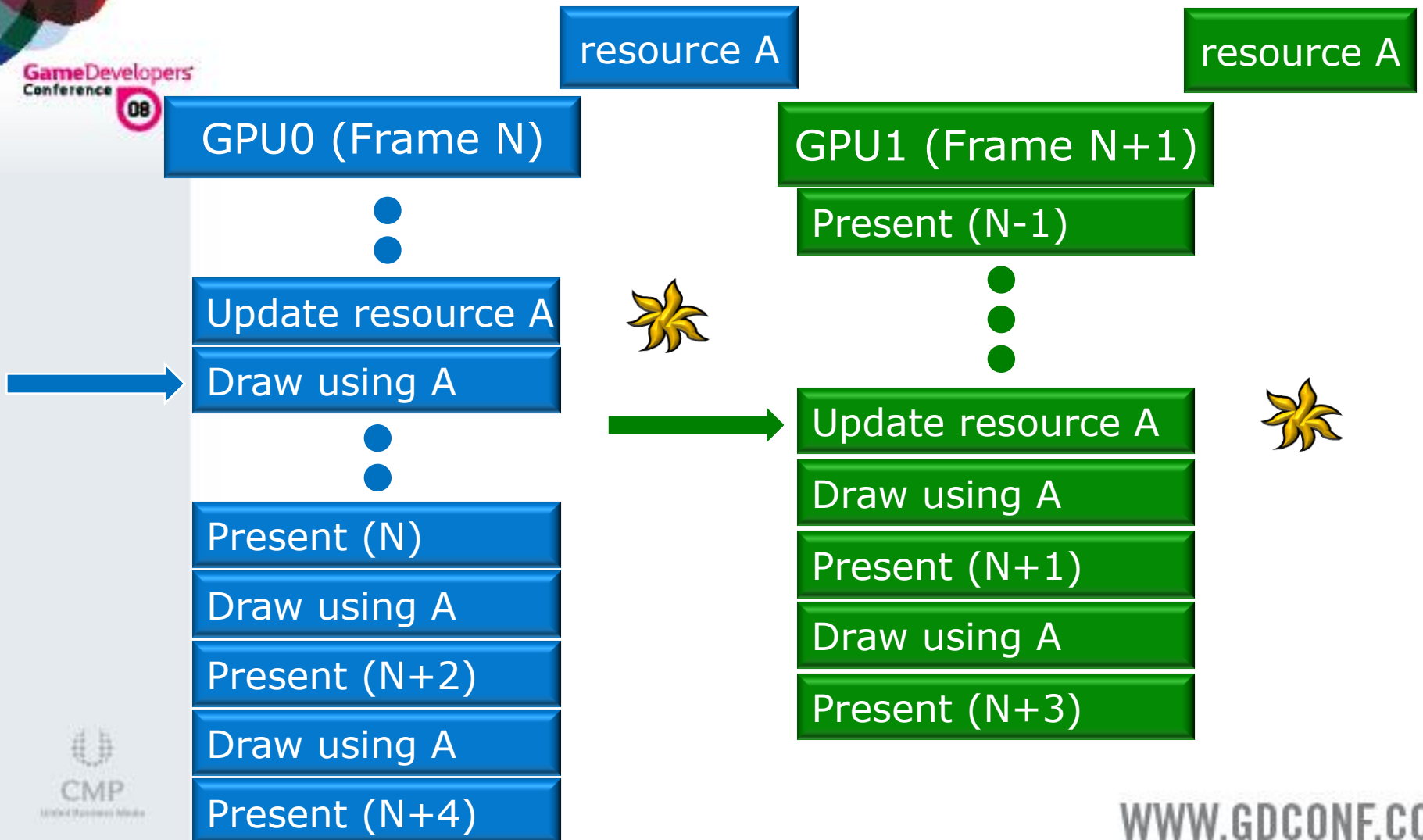
08

Solution: Resources that Change Every Few Frames



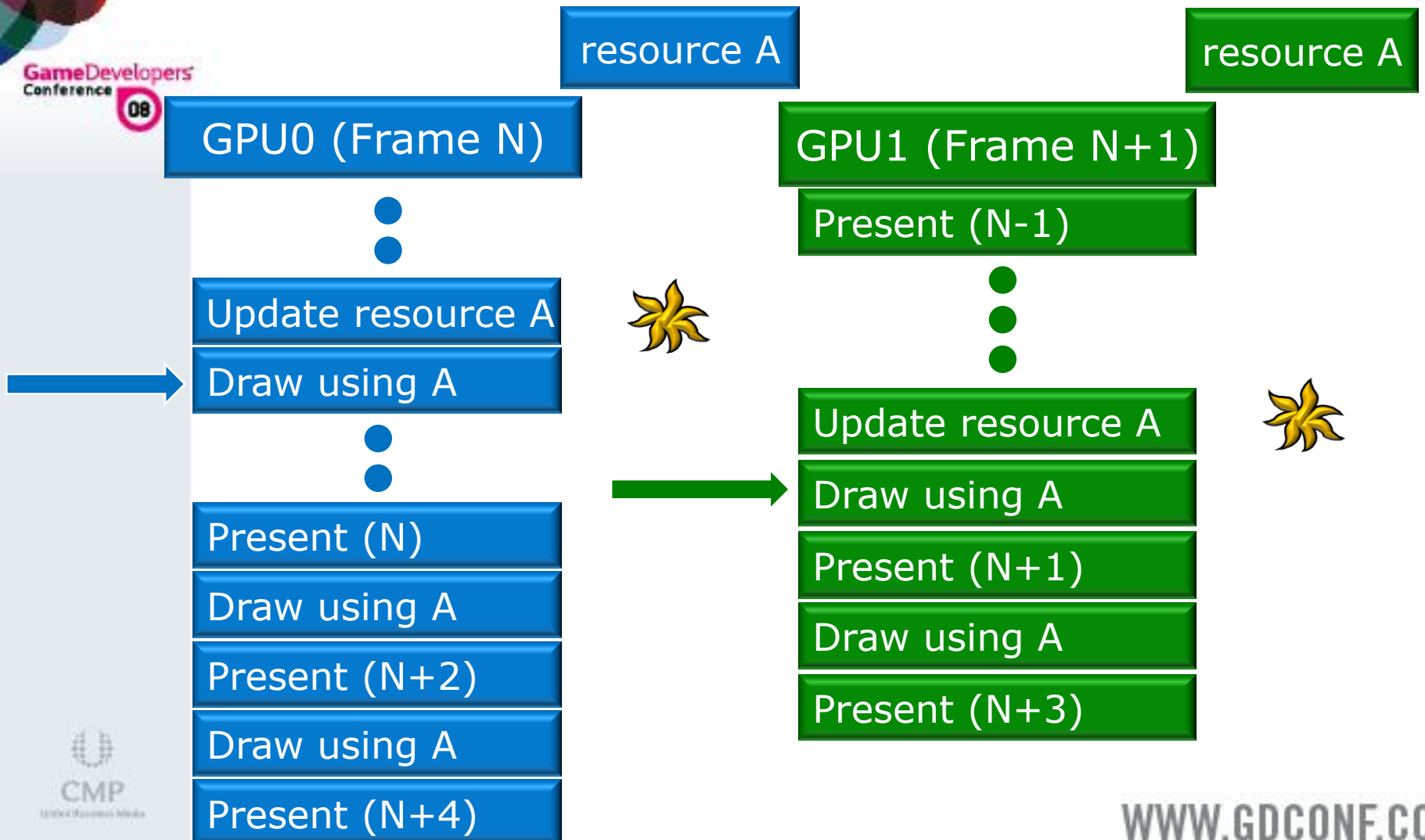


Solution: Resources that Change Every Few Frames



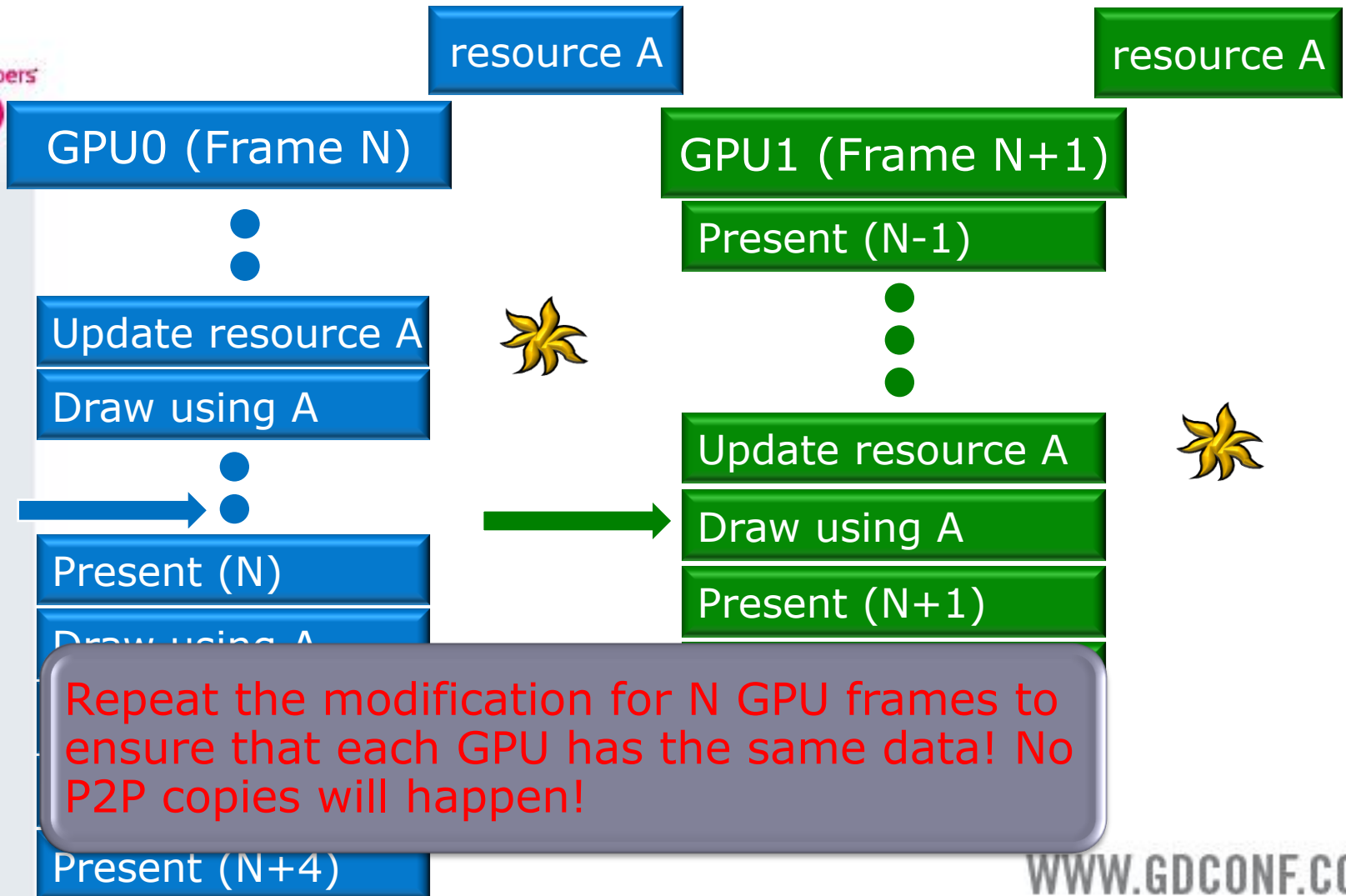


Solution: Resources that Change Every Few Frames





Solution: Resources that Change Every Few Frames





Preventing P2P copies

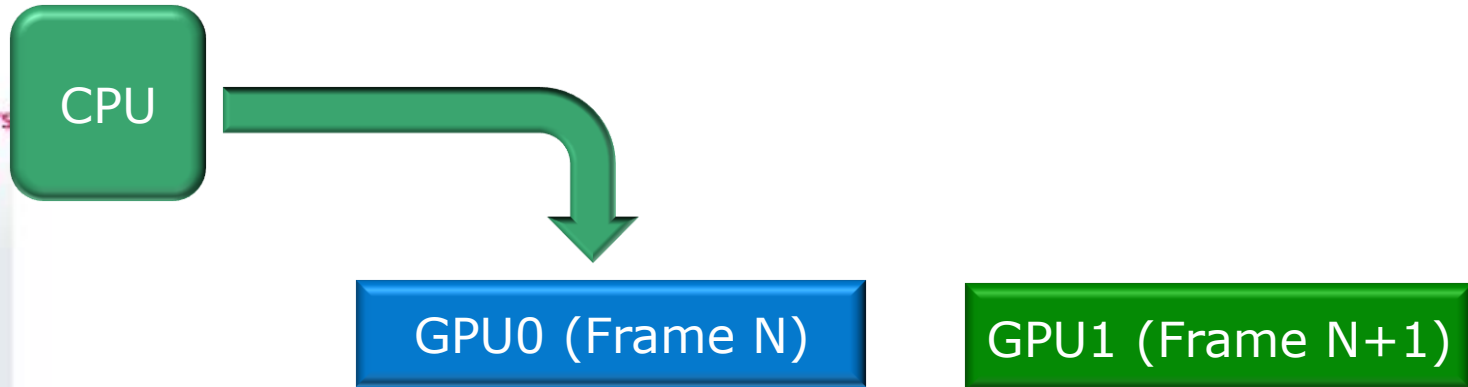
- ③ NVIDIA: Make sure to Clear RTs that don't need to be in sync
- ③ NVIDIA: Use Map(WRITE_DISCARD) to also hint resources that don't need to be in sync



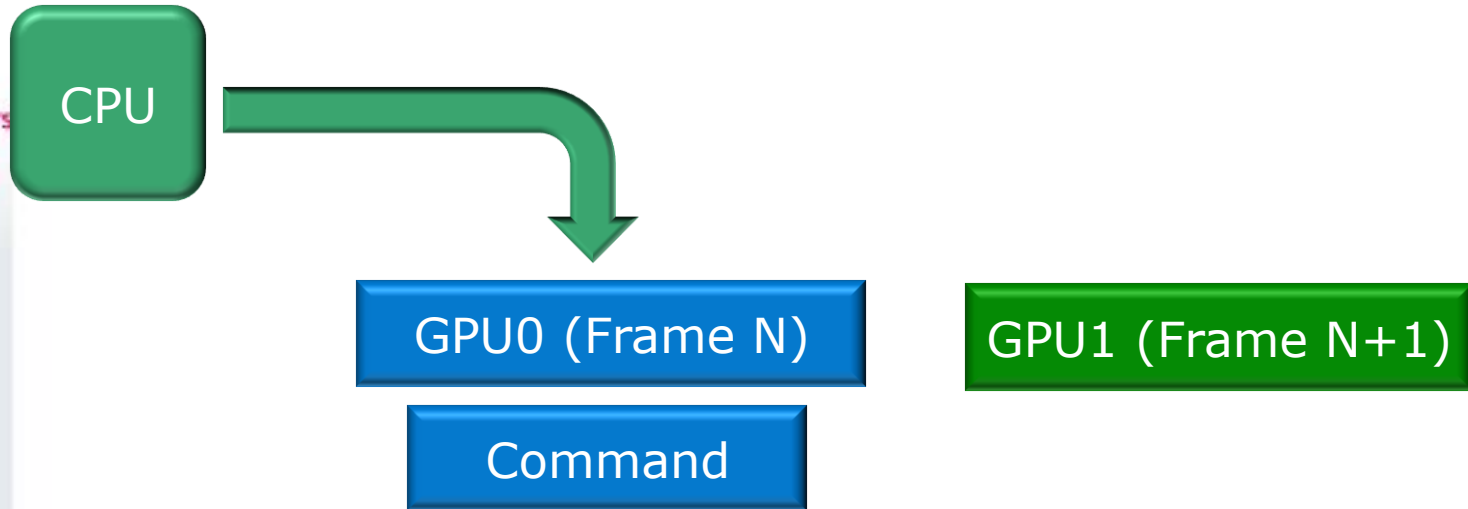
Pitfalls: In DX10 there are Other Ways to Update Resources...

- ⦿ Drawing to vertex/index buffers
- ⦿ Stream Out
- ⦿ CopyResource() calls
- ⦿ CopySubresourceRegion() calls
- ⦿ GenerateMips() calls
- ⦿ ResolveSubresource() calls
- ⦿ Do not use same resource as destination of both Map(WRITE_DISCARD) and CopyResource/CopySubresourceRegion calls

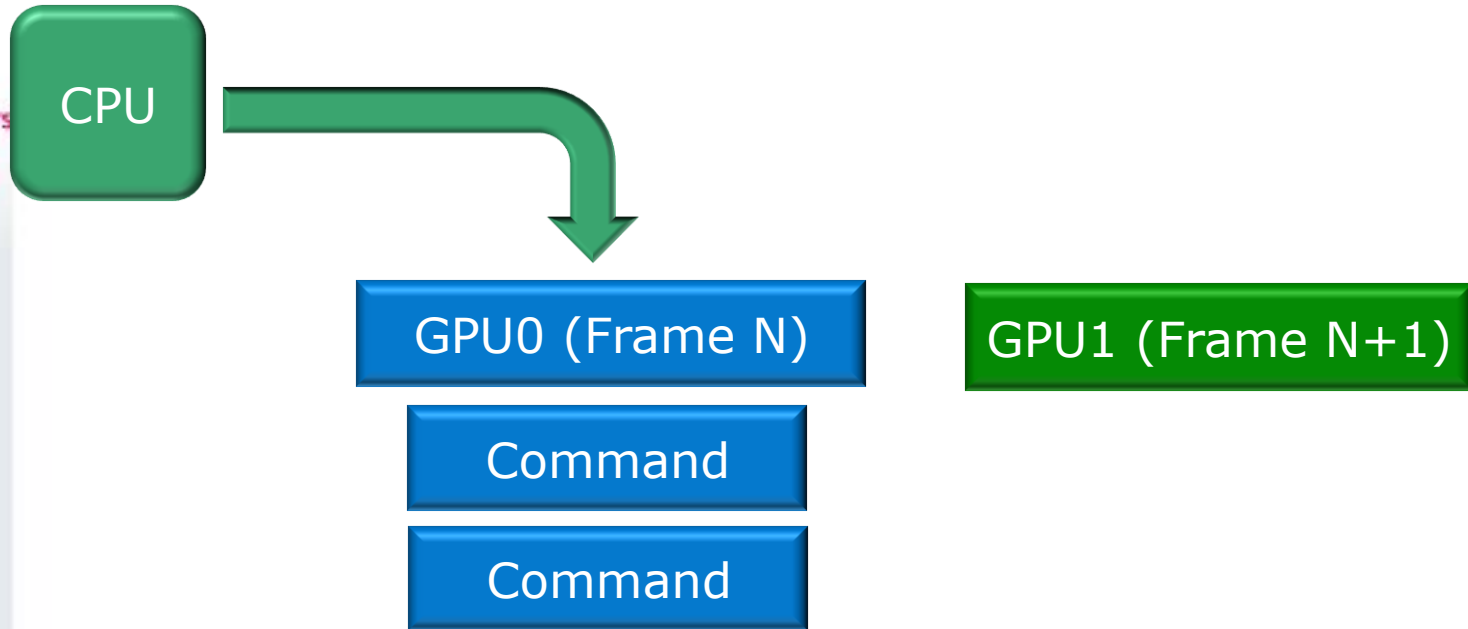
Pitfall: Waiting on Queries



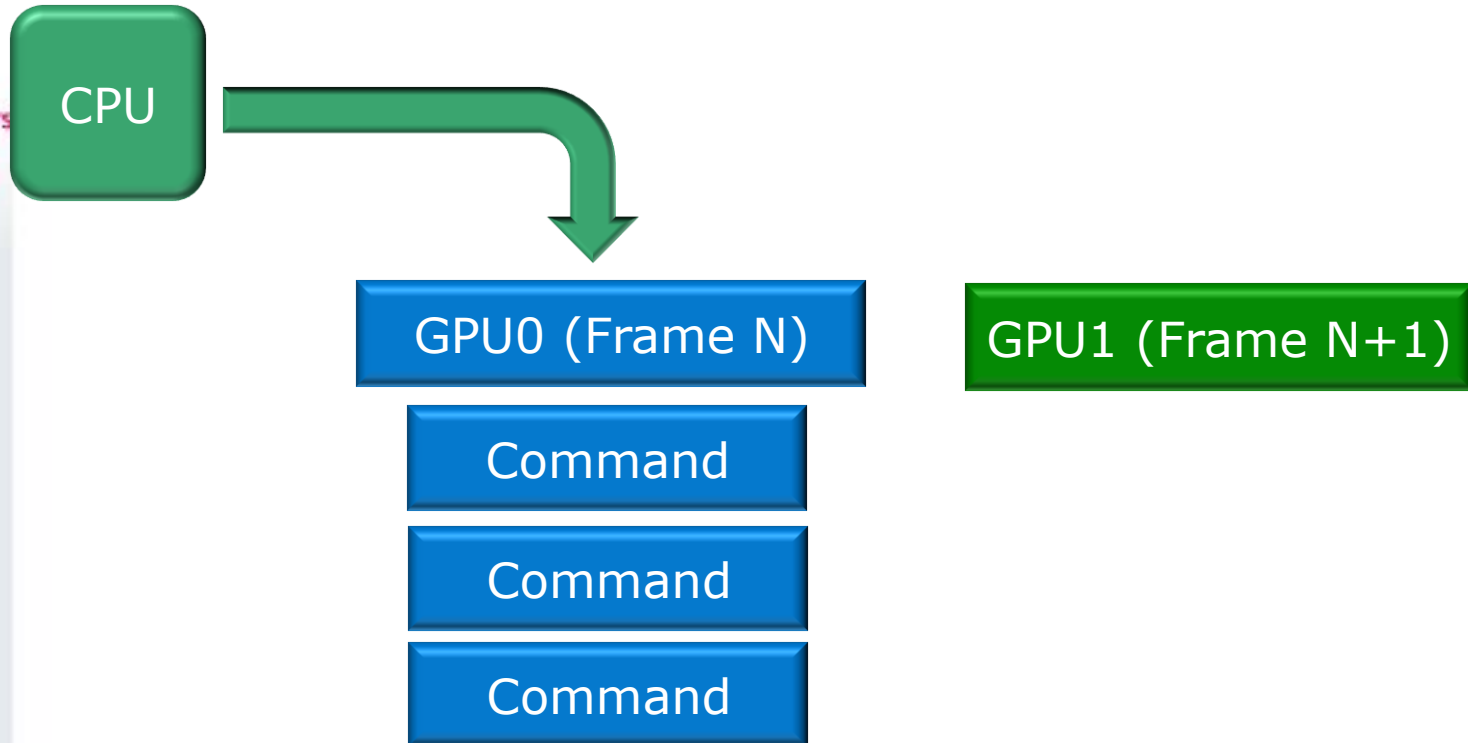
Pitfall: Waiting on Queries



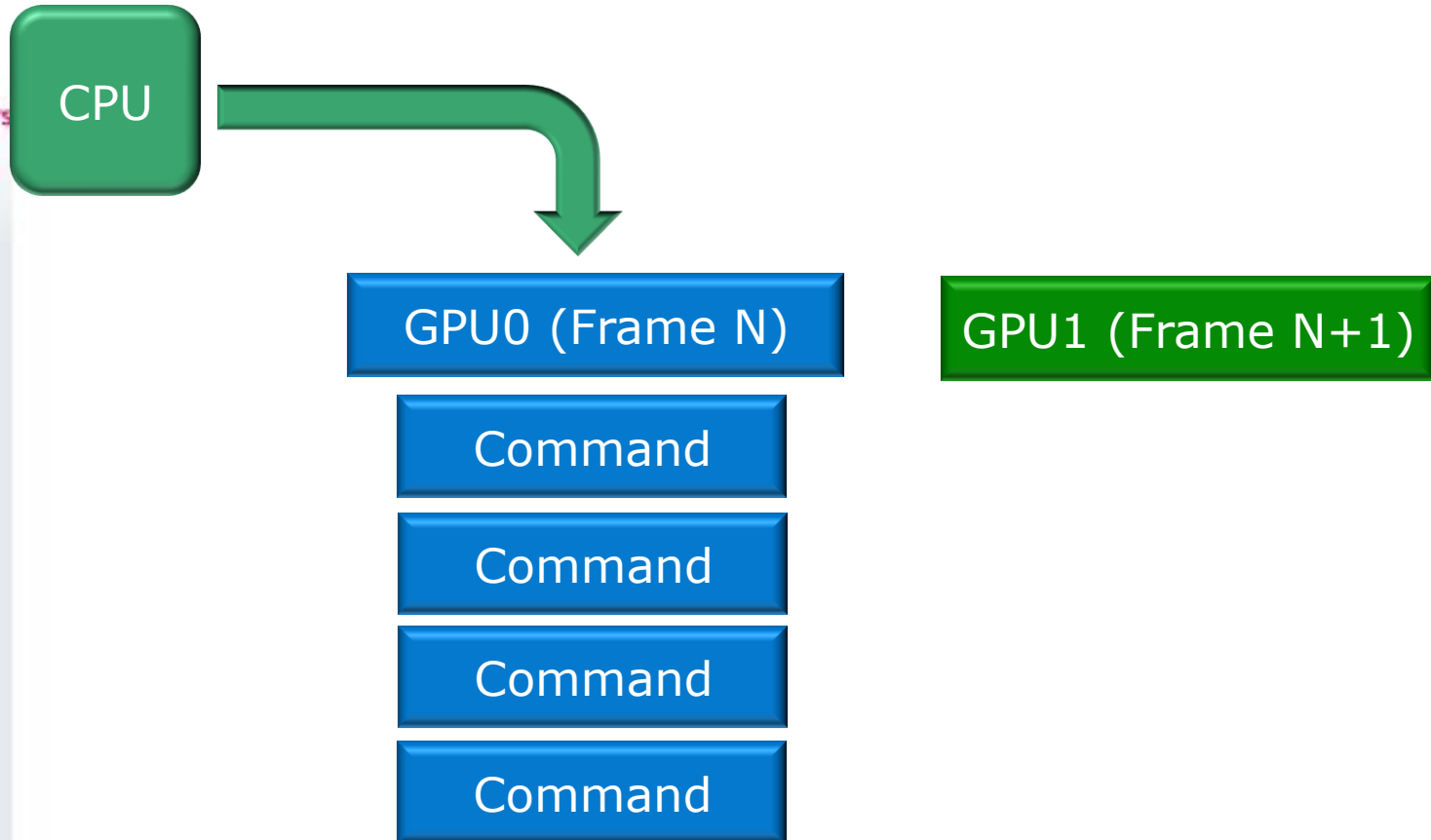
Pitfall: Waiting on Queries



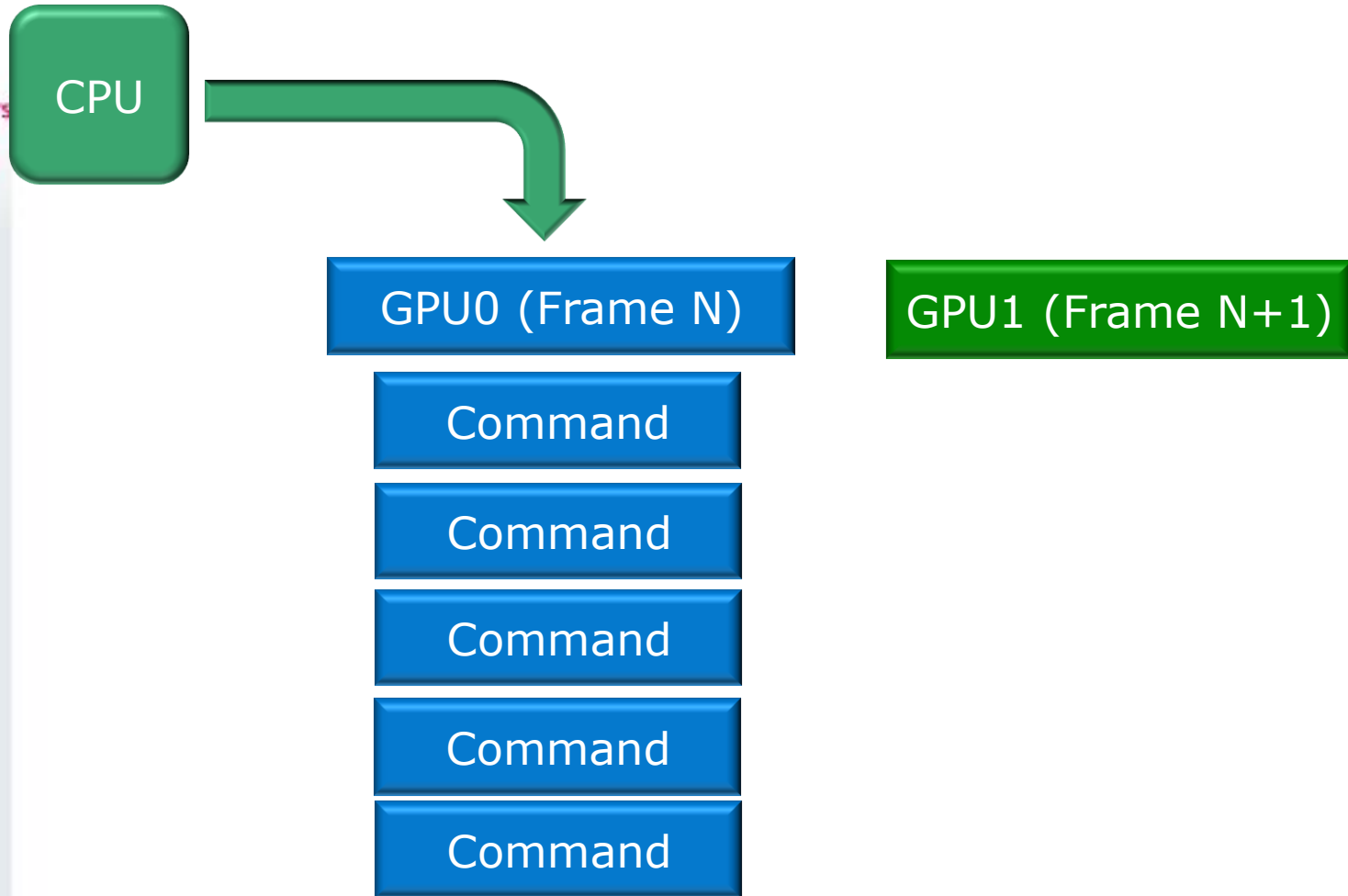
Pitfall: Waiting on Queries



Pitfall: Waiting on Queries

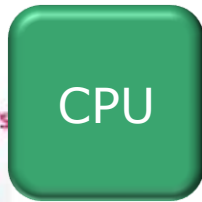


Pitfall: Waiting on Queries



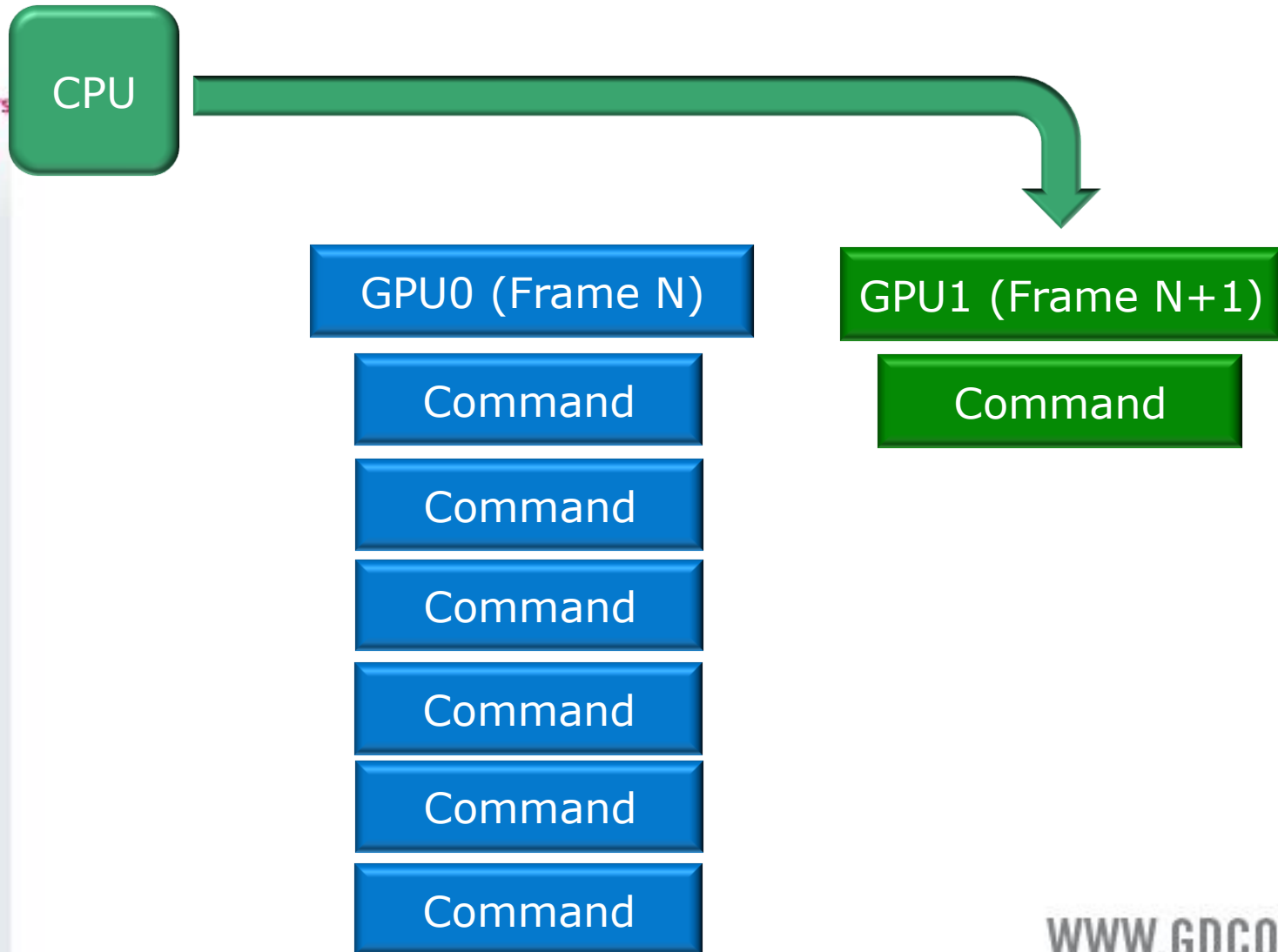


Pitfall: Waiting on Queries



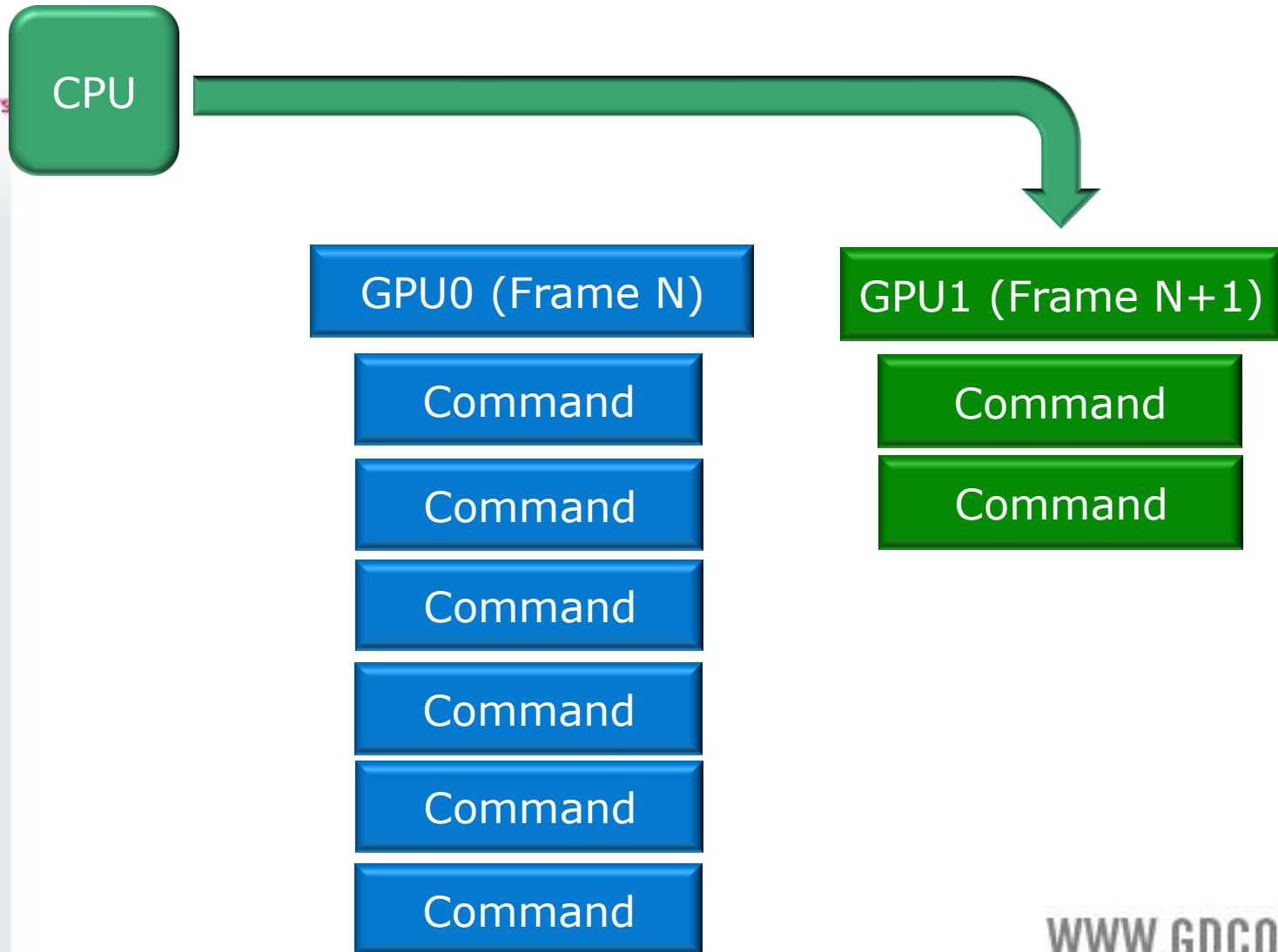


Pitfall: Waiting on Queries



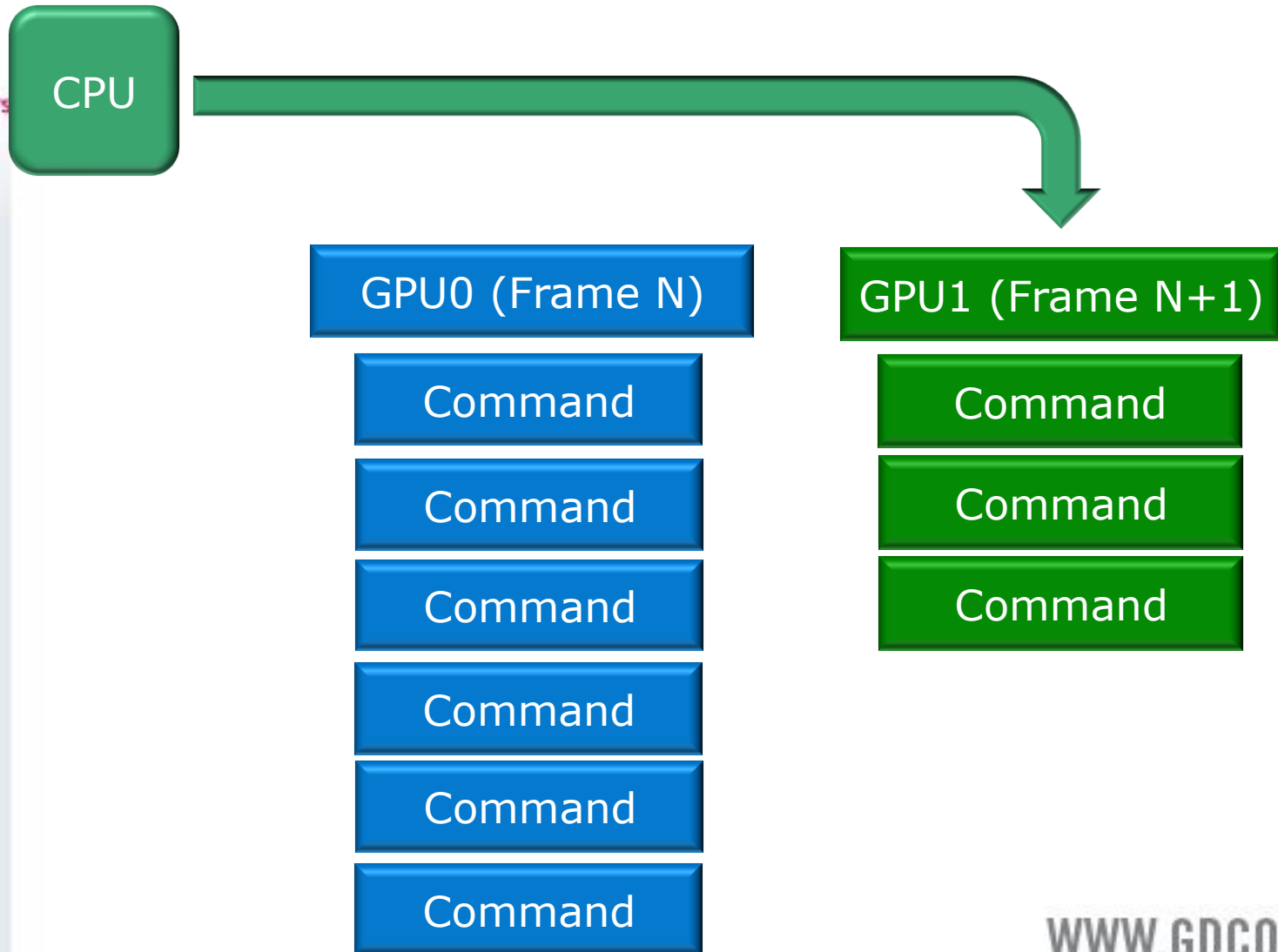


Pitfall: Waiting on Queries



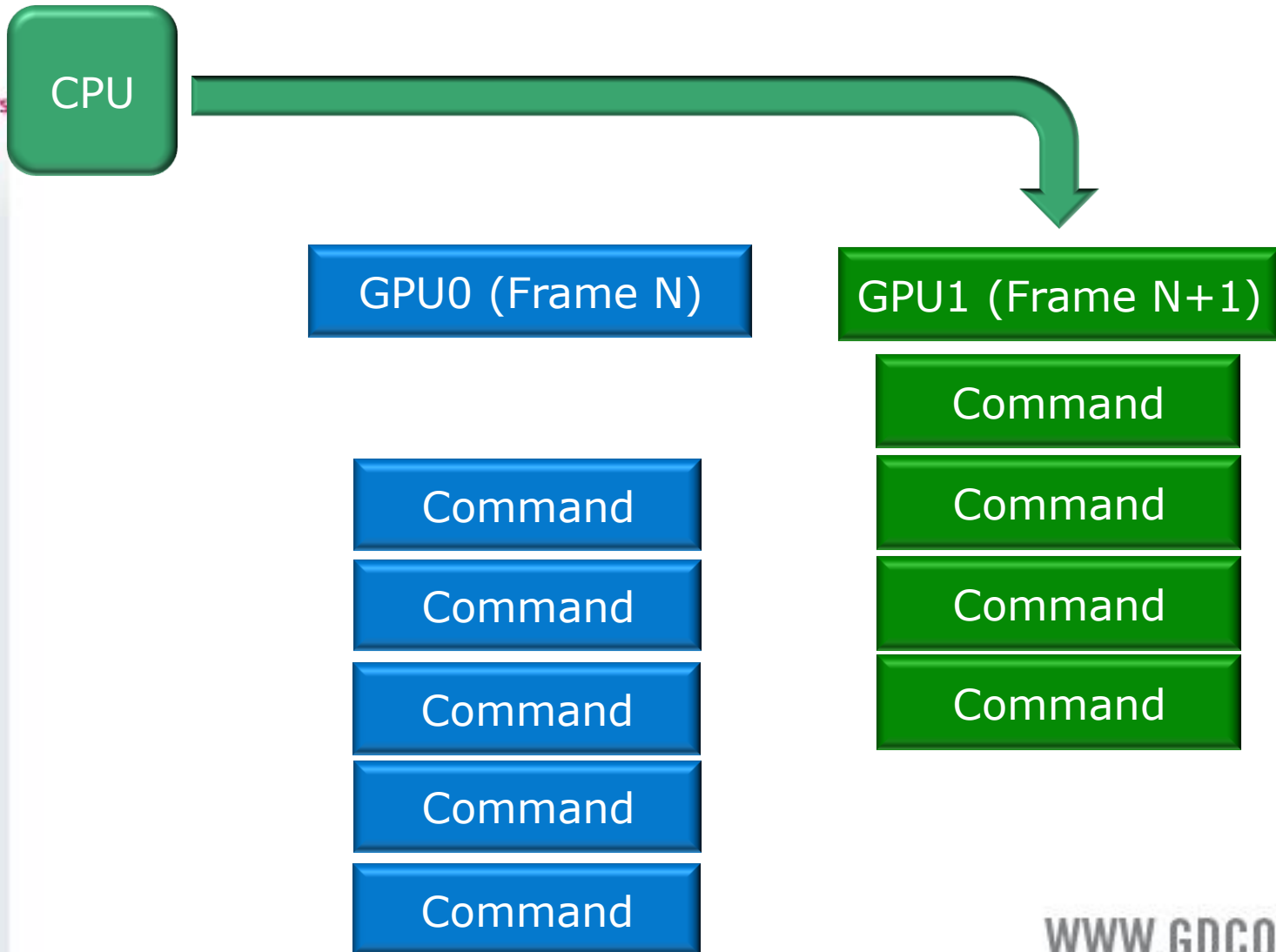


Pitfall: Waiting on Queries



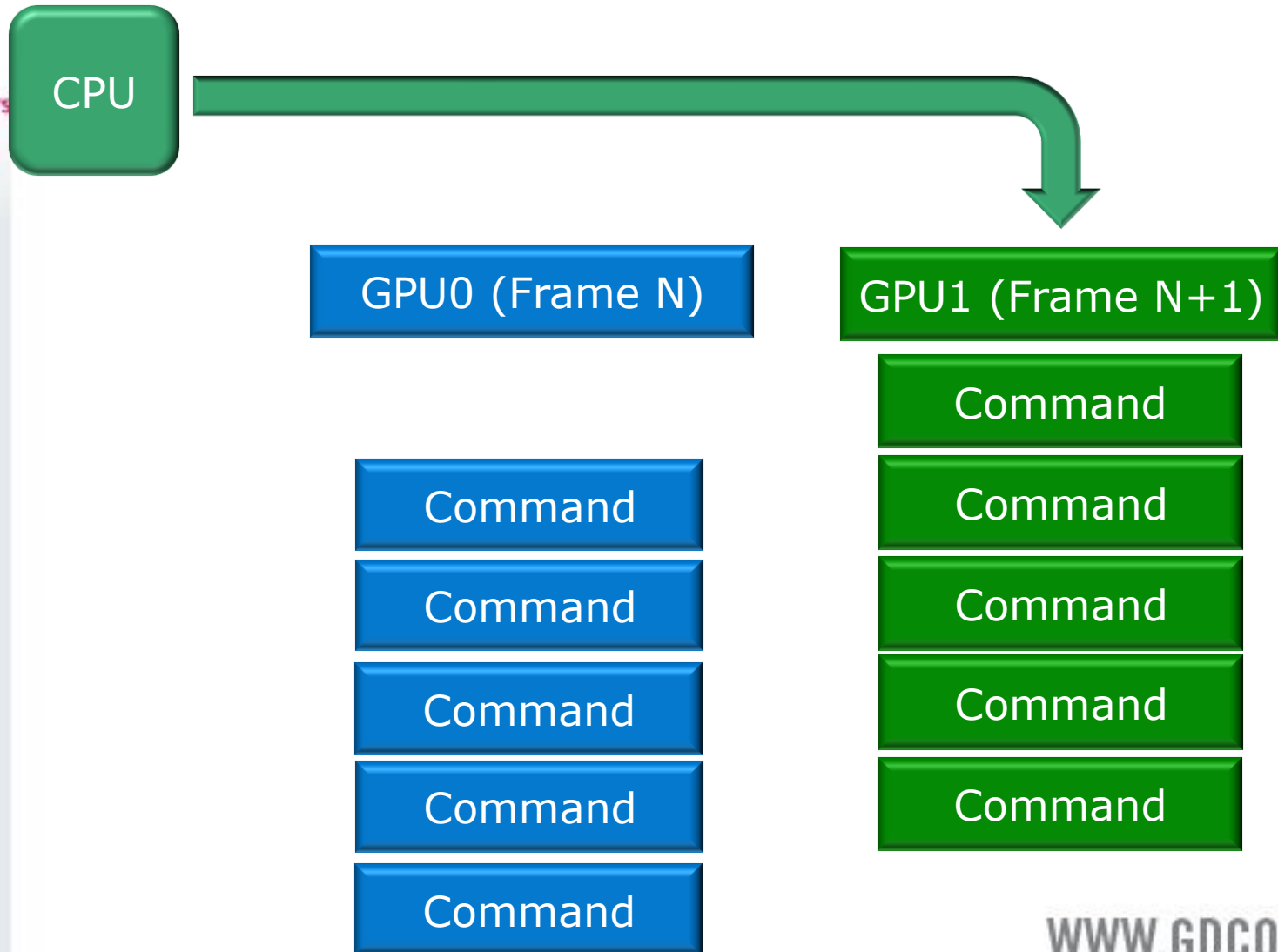


Pitfall: Waiting on Queries



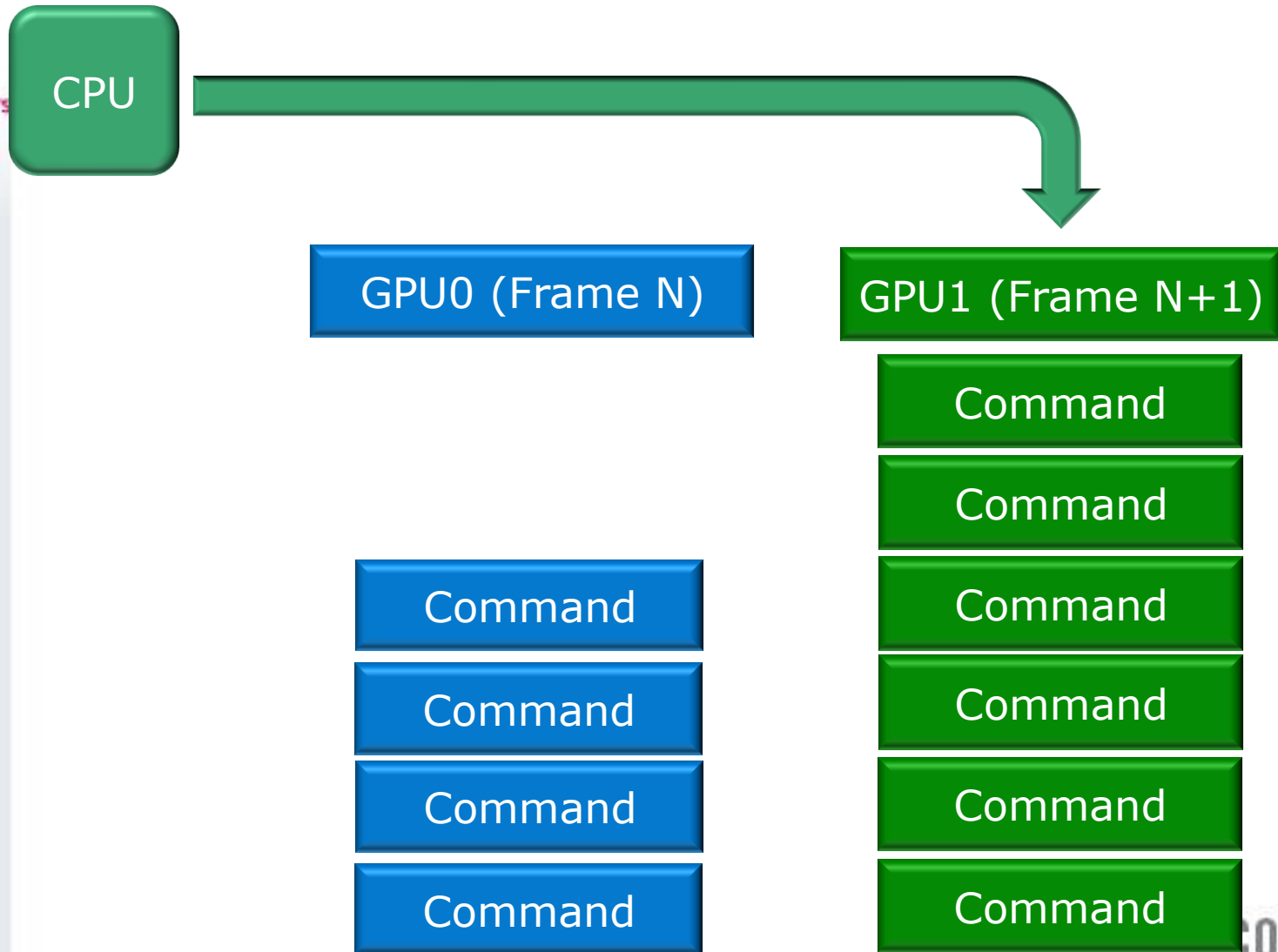


Pitfall: Waiting on Queries



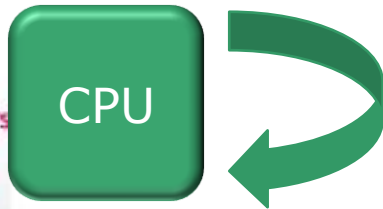


Pitfall: Waiting on Queries





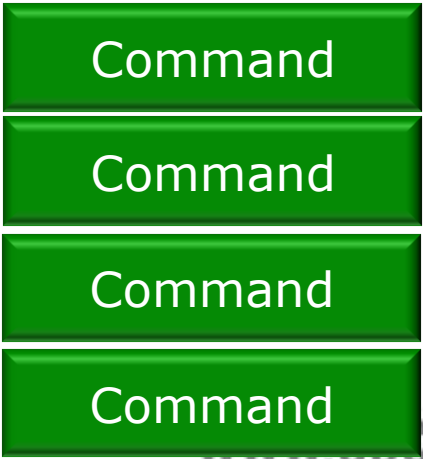
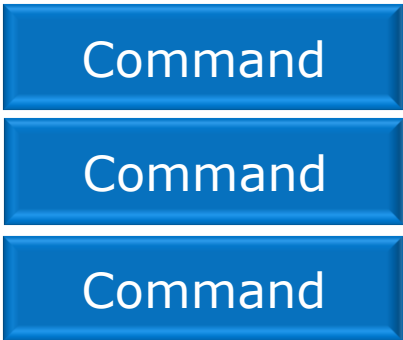
Pitfall: Waiting on Queries



Waiting for Query Result!!!

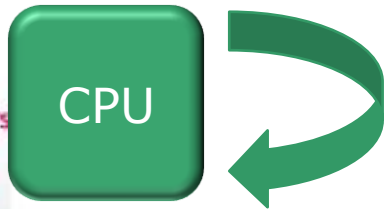
GPU0 (Frame N)

GPU1 (Frame N+1)





Pitfall: Waiting on Queries



Waiting for Query Result!!!

GPU0 (Frame N)

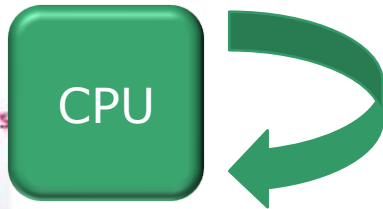
GPU1 (Frame N+1)

Command
Command

Command
Command
Command



Pitfall: Waiting on Queries



Waiting for Query Result!!!

GPU0 (Frame N)

GPU1 (Frame N+1)

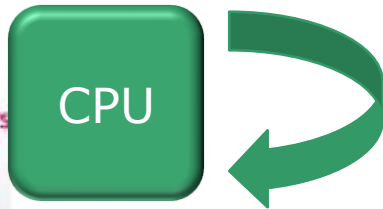
Command

Command

Command



Pitfall: Waiting on Queries



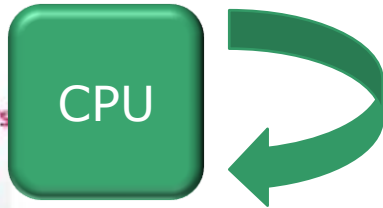
Waiting for Query Result!!!

GPU0 (Frame N)

GPU1 (Frame N+1)



Pitfall: Waiting on Queries



Waiting for Query Result!!!

GPU0 (Frame N)

GPU1 (Frame N+1)

Waiting starves GPU queues

Waiting limits parallelism

Waiting => CPU limitation



Solutions for Queries

- ③ Avoid using queries wherever possible
- ③ Make sure that a matched pair of BeginQuery / EndQuery calls occur within the same frame
- ③ Avoid waiting on query results
- ③ For queries issued every frame
 - ③ Create additional query objects for each GPU
 - ③ Cycle through them
- ③ Pickup the result of a query at least N-GPU frames after it was issued
- ③ For occlusion queries consider a CPU based approach



Pitfall: CPU Access to a Renderable Resource

- ⌚ When the CPU locks a renderable resource it must wait for all GPUs to finish using the resource before acquiring the pointer
- ⌚ All GPUs now have to wait until the CPU unlocks the resource pointer
- ⌚ After the unlock the driver has to update the resource on each GPU via P2P copies
- ⌚ Just don't do this – it destroys performance even on a single GPU setup, and is catastrophic for MGPUs



Solutions: Locks / Maps

- ④ In DX10 stream to and copy from STAGING textures
- ④ In DX9 StretchRect() is always better than Lock()
- ④ At resource creation time use the appropriate flags from:
 - ④ D3D10_USAGE
 - ④ D3D10_CPU_ACCESS_FLAG
- ④ **AMD: In DX9 never lock static Vertex/Index Buffers because it will cause P2P copies**



Concluding Pitfalls & Solutions

- ③ Drivers take a conservative approach
 - ③ Performing P2P copies
- ③ You know the application best
 - ③ Determine if a P2P copy is necessary
 - ③ Talk to us about an app profile
- ③ Never use resources shared between devices on DX10
 - ③ No way to detect update by other app



Call to Action

- ④ MGPUs provide demonstrable performance/quality gains even on low end setups
- ④ Plan from day one to make your rendering scale
 - ④ Use AMD library to detect number of GPUs
 - ④ Use NVIDIA library to detect MGPU topology
- ④ Regularly check for AFR unfriendly behavior
- ④ Talk to us holger.gruen@amd.com
illamas@nvidia.com



Please Fill In The Feedback Forms



Questions?